Name: **Sachin Mishra**
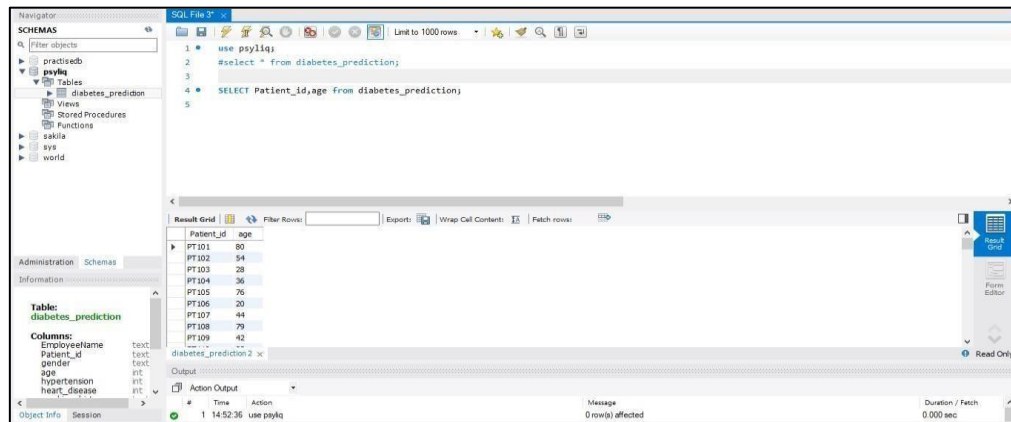Email: smzp274@gmail.com
LinkedIn: https://www.linkedin.com/in/sachin-mishra-4169b7213/
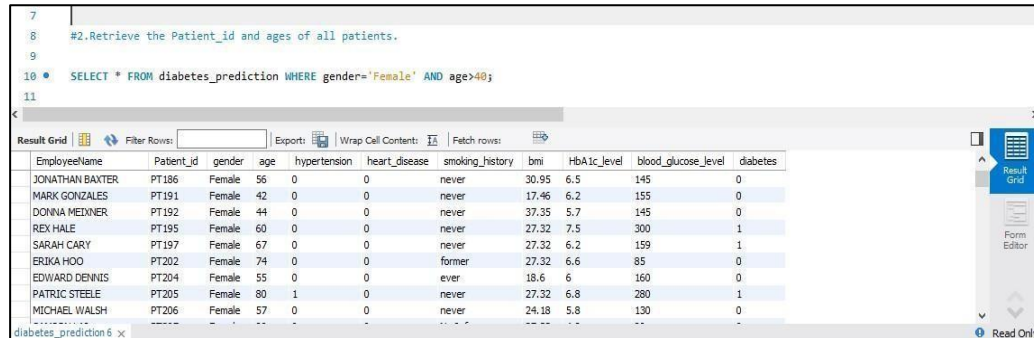
# DIABETES  PREDICTION  ASSESSMENT – Psyliq

1. Retrieve the Patient_id and ages of all patients.

   SELECT Patiend_id, age FROM diabetes_prediction;



2. Select all female patients who are older than 40.
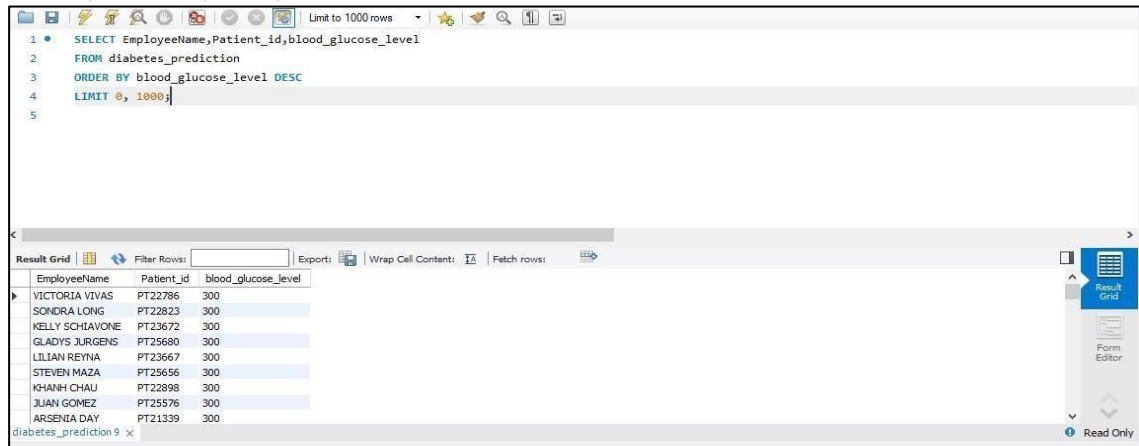   SELECT * FROM diabetes_prediction WHERE gender='Female' AND age>40;



3. Calculate the average BMI of patients.
   SELECT
   EmployeeName,Patient_id,AVG(bmi)
   FROM diabetes_prediction
   GROUP BY
   EMPLOYEENAME,Patiend_id
   LIMIT 0,1000;

4. List patients in descending order of blood glucose levels.

SELECT EmployeeName, Patient_id,
blood_glucose_levelFROM diabetes_prediction
ORDER BY blood_glucose_level
DESC,LIMIT0,1000;

```
1 •   SELECT EmployeeName,Patient_id,blood_glucose_level
2     FROM diabetes_prediction
3     ORDER BY blood_glucose_level DESC
4     LIMIT 0, 1000;
5
```
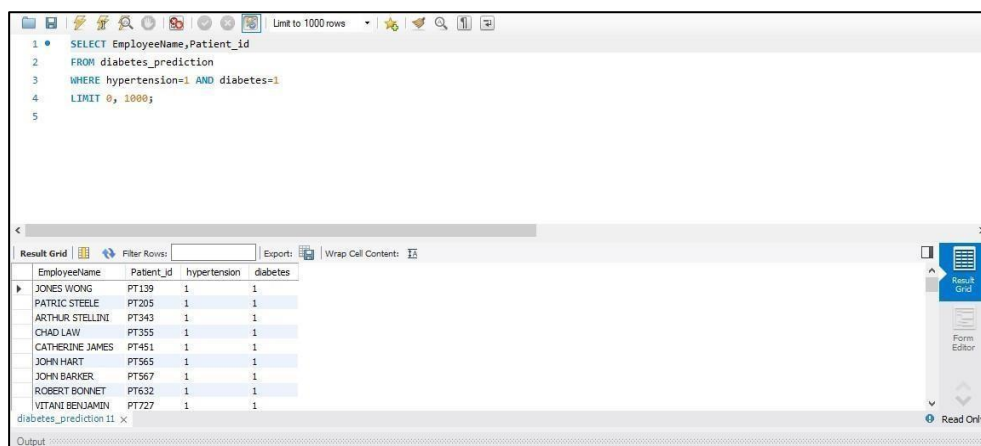
Result Grid | Filter Rows: | Export: | Wrap Cell Content: | Fetch rows:

| EmployeeName | Patient_id | blood_glucose_level |
|---|---|---|
| VICTORIA VIVAS | PT22786 | 300 |
| SONDRA LONG | PT22823 | 300 |
| KELLY SCHIAVONE | PT23672 | 300 |
| GLADYS JURGENS | PT25680 | 300 |
| LILIAN REYNA | PT23667 | 300 |
| STEVEN MAZA | PT25656 | 300 |
| KHANH CHAU | PT22898 | 300 |
| JUAN GOMEZ | PT25576 | 300 |
| ARSENIA DAY | PT21339 | 300 |

diabetes_prediction 9 ×

5. Find patients who have hypertension and diabetes.

SELECT
EmployeeName,Patient_id
FROM diabetes_prediction
WHERE hypertension=1 AND
diabetes=1LIMIT0,1000;

```
1 •   SELECT EmployeeName,Patient_id
2     FROM diabetes_prediction
3     WHERE hypertension=1 AND diabetes=1
4     LIMIT 0, 1000;
5
```
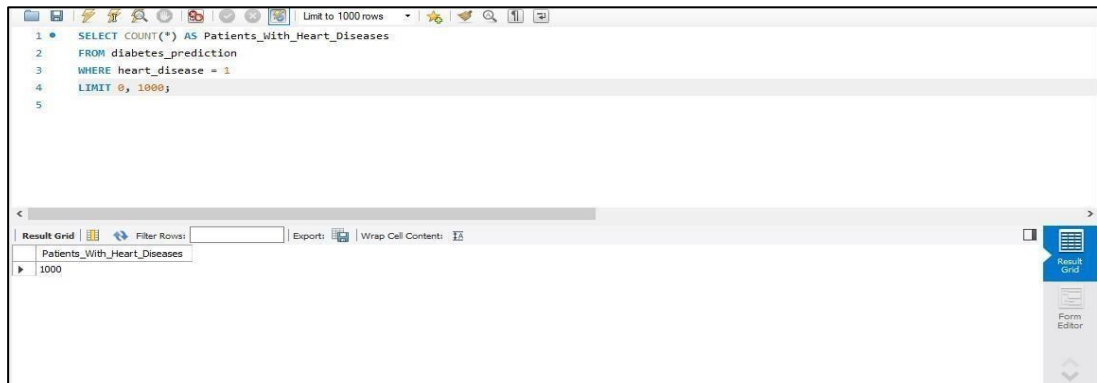
Result Grid | Filter Rows: | Export: | Wrap Cell Content:

| EmployeeName | Patient_id | hypertension | diabetes |
|---|---|---|---|
| JONES WONG | PT139 | 1 | 1 |
| PATRIC STEELE | PT205 | 1 | 1 |
| ARTHUR STELLINI | PT343 | 1 | 1 |
| CHAD LAW | PT355 | 1 | 1 |
| CATHERINE JAMES | PT451 | 1 | 1 |
| JOHN HART | PT565 | 1 | 1 |
| JOHN BARKER | PT567 | 1 | 1 |
| ROBERT BONNET | PT632 | 1 | 1 |
| VITANI BENJAMIN | PT727 | 1 | 1 |

diabetes_prediction 11 ×

Output

6. Determine the number of patients with heart
   disease.
   SELECT COUNT(*) AS
    Patients_With_Heart_diseasesFROM
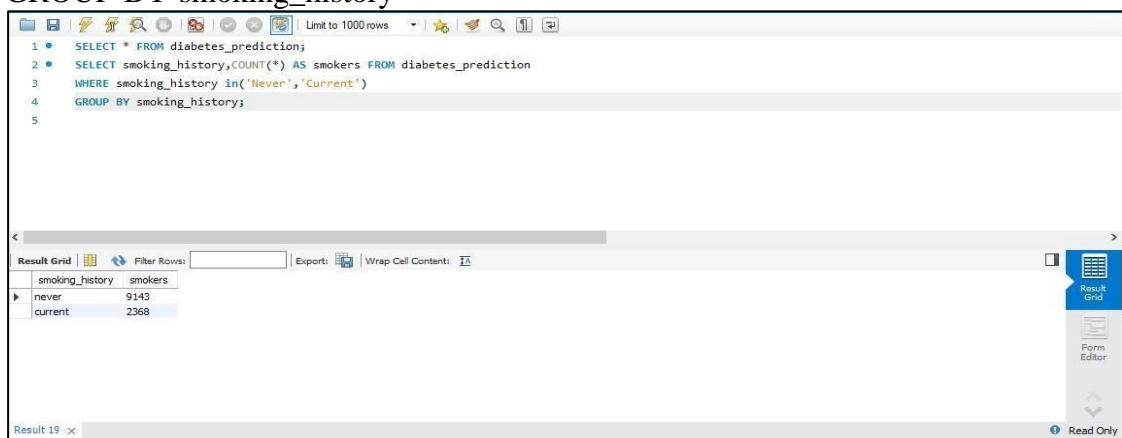    diabetes_prediction
   WHERE heart_disease=1

7. Group patients by smoking history and count how many smokers and nonsmokers there are.

SELECT * FROM diabetes_prediction;
SELECT smoking_history, COUNT(*) AS smokers
FROM diabetes_prediction
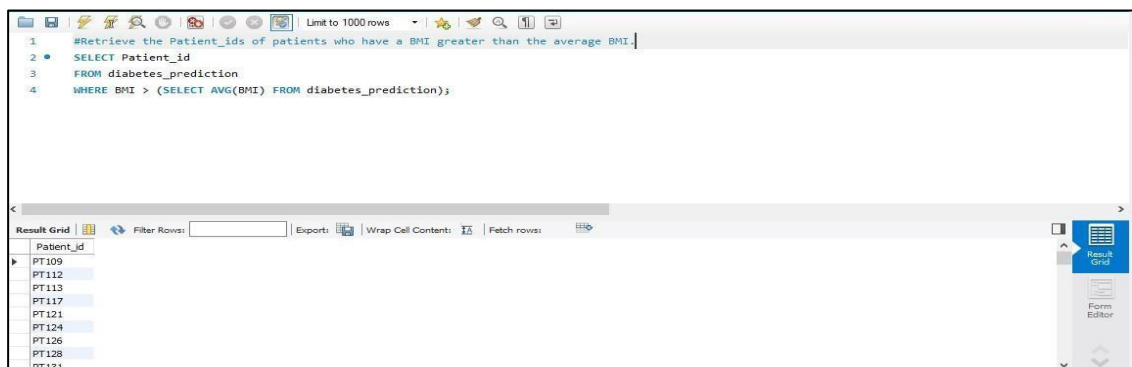WHERE smoking_history in("Never","Current")
GROUP BY smoking_history



8. Retrieve the Patient_ids of patients who have a BMI greater than the average BMI.

SELECT Patient_id

FROM diabetes_prediction
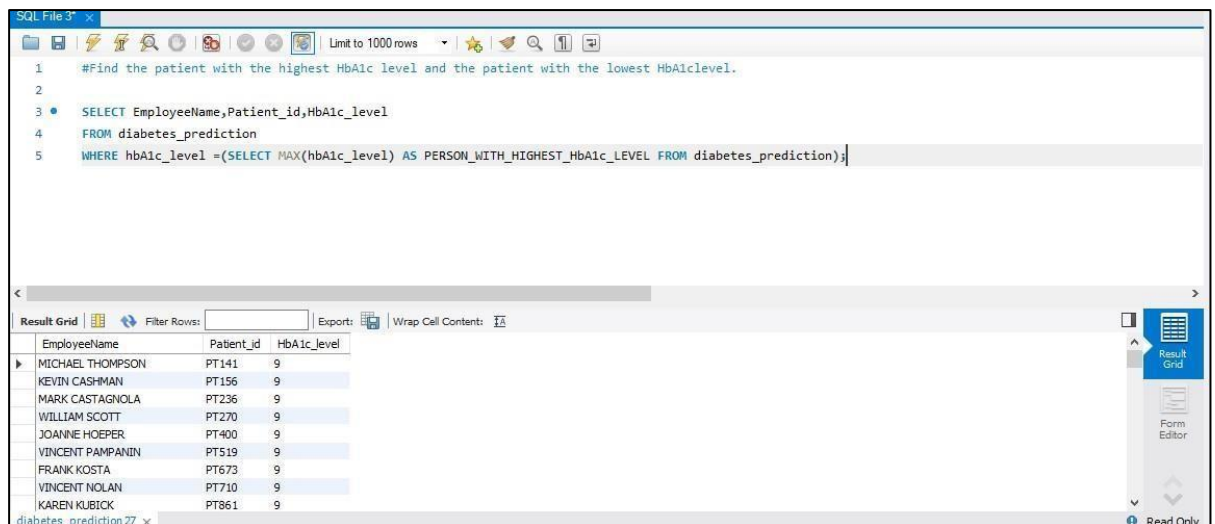
WHERE BMI > (SELECT AVG(BMI) FROM diabetes_prediction);

9. Find the patient with the highest HbA1c level and the patient with the lowest HbA1clevel.

## Highest hbA1c_level

SELECT
EmployeeName,Patient_id,HbA1c_levelFROM
diabetes_prediction
WHERE
hbA1c_level=(SELECTMAX(hbA1c_level)AS
PERSON_WITH_HIGHEST_hbA1c_LEVEL
FROM diabetes_prediction;



## Lowest hbA1c_level

SELECT
EmployeeName,Patient_id,HbA1c_level FROM
diabetes_prediction
WHEREhbA1c_level=(SELECTMIN(hbA1c_lev
el)AS
PERSON_WITH_LOWEST_hbA1c_LEVEL
FROM diabetes_prediction;

10. Calculate the age of patients in years (assuming the current date as of now).
    SELECT Patient_id,
    DATE_FORMAT(NOW(),'%Y') – DATE_FORMAT(age,'%Y') –
    DATE_FORMAT(NOW(),'00-%m-
    %d')   <        DATE_FORMAT(age, '00-%m-%d'))AS Age
    FROM diabetes_prediction



11. Rank patients by blood glucose level within each gender group.

    SELECT EmployeeName,Patient_id,gender,blood_glucose_level,
    RANK() OVER(PARTITION BY gender ORDER BY blood_glucose_level
    DESC)ASRank_ By_Blood_Glucose_Level;

12. Update the smoking history of patients who are older than 50 to "Ex-smoker."

```
SELECT * FROM diabetes_prediction;UPDATE
diabetes_prediction
SET smoking_history = "EX-
Smoker"WHEREage>50;
```



13. Insert a new patient into the database with sample data.

```
INSERT INTO diabetes_prediction (EmployeeName, Patient_id,
gender,age, hypertension,heart_disease, smoking_history, bmi,
 HbA1c_level, blood_glucose_level, diabetes)
VALUES ('Peter Hartley','PT100101','Male',45,1,0,'Never',12.34,5,136,0);
```

14. Delete all patients with heart disease from the database.

```
DELETE FROM
diabetes_predictionWHERE
heart_disease=1;
```



15. Find patients who have hypertension but not diabetes using the EXCEPT operator.

```
SELECT * FROM
diabetes_predictionWHERE
hypertension=1
EXCEPT
SELECT * FROM diabetes_prediction WHERE diabetes=1;
```

16. Define a unique constraint on the "patient_id" column to ensure its values
    are unique.
    ALTER TABLE diabetes_prediction
    MODIFY COLUMN Patient_id VARCHAR(100);

    ALTER TABLE diabetes_prediction
    ADD CONSTRAINT UPID UNIQUE(Patient_id);



17. Create a view that displays the Patient_ids, ages, and BMI of patients.

    CREATE view PatientInfo AS SELECT
    Patient_id,age,bmiFROM diabetes_prediction;

    SELECT * from PatientInfo;

18.     Suggest improvements in the database schema to reduce data redundancy and improve data integrity.

Reducing data redundancy and enhancing data integrity in a database schema involves implementing the following improvements:

1.  Normalization:
   Apply normalization techniques to organize data efficiently. Break down large tables into smaller ones and establish relationships to minimize redundancy. This helps in maintaining data integrity by ensuring that each piece of information is stored in only one place.

2.  Use of Primary and Foreign Keys:
   Define and utilize primary keys to uniquely identify records within a table. Implement foreign keys to establish relationships between tables, ensuring accurate linkages and preventing orphaned records.

3.  Elimination of Transitive Dependencies:
   Identify and eliminate transitive dependencies, where data is indirectly dependent on other data. This helps in removing unnecessary redundancy and inconsistencies in the database, leading to improved data integrity.

4.  Views:
   Create views to present a consolidated and normalized perspective of the data. Views act as virtual tables, reducing redundancy by providing a unified source of truth for related information and simplifying data access.

5.  Referential Integrity:
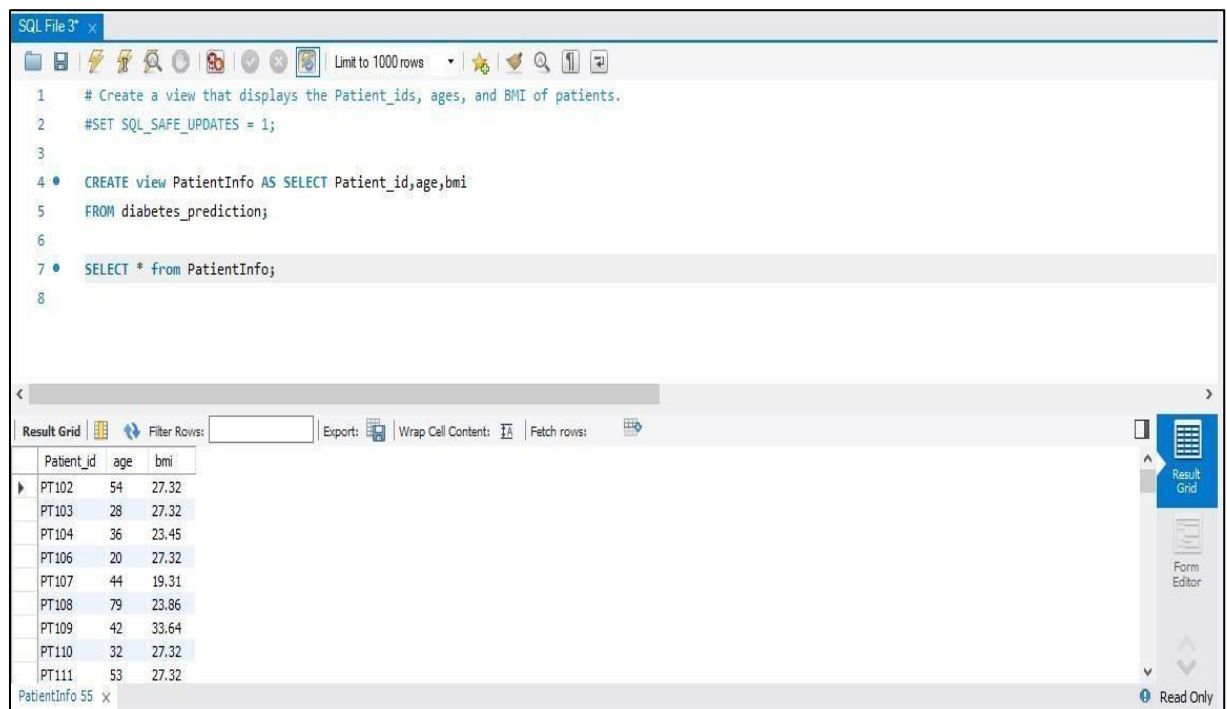   Enforce referential integrity by ensuring that foreign keys reference valid primary keys. This prevents inconsistencies, orphaned records, and strengthens the overall integrity of the data.

6.  Triggers:
   Implement triggers to automate actions based on specified events, such as inserting or updating records. Triggers help enforce business rules and maintain data integrity by automatically updating related information, reducing redundancy.

7.  Indexes:
   Utilize indexes to optimize query performance. Indexes allow for quick data retrieval based on specified criteria, reducing the need to scan through large datasets and minimizing redundancy in data access.

8.  Data Validation:
   Implement robust data validation rules to ensure the accuracy and consistency of data entered into the database. By enforcing validation criteria, redundancy is minimized, reducing the likelihood of duplicate records or errors in data input.

9.  Documentation:
   Maintain comprehensive documentation of the database schema, including relationships, constraints, and considerations for developers and users. Clear documentation aids in understanding the database structure, reducing the risk of errors and enhancing data integrity.

10.  Regular Review and Update:
    Conduct periodic reviews of the database schema to identify opportunities for optimization and improvement. Adapt the schema based on evolving business requirements, ensuring that it continues to support data integrity goals over time.

Incorporating these improvements significantly reduces data redundancy, enhances data integrity, and establishes a more efficient and reliable database schema.

**19. Explain how you can optimize the performance of SQL queries on this dataset.**

Optimizing the performance of SQL queries on a dataset involves several strategies to ensure efficient data retrieval and processing. Here are some key approaches:

1.  Indexing:
    - Identify columns frequently used in WHERE clauses, JOIN conditions, or ORDER BY clauses.
    - Create indexes on these columns to speed up data retrieval.
    - Be cautious not to over-index, as it may impact write performance.

2.  Query Rewriting:
    - Review and rewrite complex queries to simplify their structure.
    - Use appropriate JOIN types and conditions to minimize the amount of data processed.
    - Consider breaking down complex queries into simpler, well-optimized subqueries.

3.  Avoid SELECT *:
    - Instead of selecting all columns using `SELECT *`, explicitly list only the necessary columns.
    - This reduces the amount of data transferred and improves query performance.

4.  Use Proper Joins:
    - Choose the most efficient JOIN type (INNER JOIN, LEFT JOIN, etc.) based on the relationships between tables.
    - Ensure that join conditions are properly indexed.

5.  Aggregate Functions:
    - Minimize the use of expensive aggregate functions like COUNT DISTINCT or GROUP BY when not essential.
    - Optimize queries that involve aggregation to reduce processing time.

6.  Partitioning:
    - If applicable, consider partitioning large tables based on certain criteria (e.g., date ranges).
    - Partition pruning can significantly reduce the amount of data scanned during queries.

7.  Update Statistics:
    - Keep database statistics up-to-date to help the query planner make informed decisions.
    - Regularly analyze and update index statistics for accurate query optimization.

8.  Limit Data Retrieval:
    - Use the LIMIT or TOP clause to restrict the number of rows returned, especially when displaying data in a user interface.
    - This reduces the load on both the database server and the network.

9.  Caching:
    - Implement caching mechanisms for frequently executed and static queries.
    - Cached results can be served quickly without hitting the database, enhancing performance.

10. Materialized Views:
    - Consider using materialized views for precomputed aggregations or complex queries.
    - Materialized views store results physically, reducing the need to recalculate values during each query.

11. Database Sharding:
   - If the dataset is massive, consider horizontal partitioning (sharding) to distribute the data across multiple servers.
   - This can improve parallel processing and reduce the load on individual nodes.

12. Regular Database Maintenance:
   - Perform routine database maintenance tasks such as vacuuming, updating statistics, and reorganizing indexes to keep the database in optimal condition.

13. Use Stored Procedures:
   - Store frequently used queries as stored procedures, which can be compiled and optimized by the database engine.
   - This reduces the overhead of query parsing and optimization.

14. Network Optimization:
   - Optimize the network infrastructure between the application server and the database server to minimize latency.

15. Hardware Scaling:
   - Consider scaling the hardware, such as adding more RAM, optimizing disk I/O, or upgrading the CPU, to improve overall system performance.

Employing a combination of these strategies significantly enhances the performance of SQL queries on your dataset. To achieve optimal results, analyze the specific characteristics of your dataset and workload, tailoring these optimizations accordingly.