**PSYLIQ**
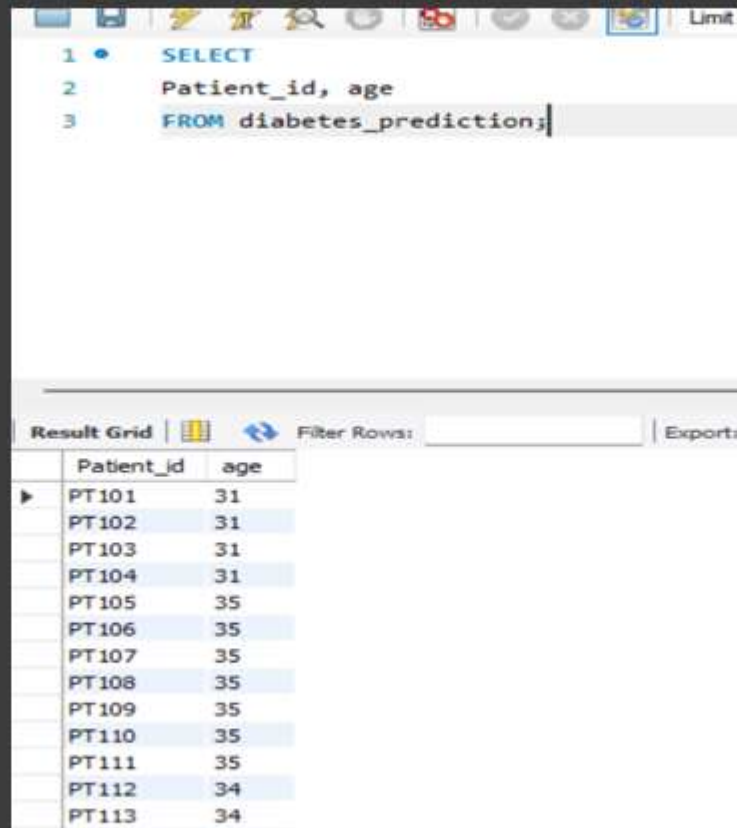
Data Analyst Internship

# TASK 3: DIABETES PREDICTION ANALYSIS
# BY
# AMAN SHAIKH

# Q1. Retrieve the Patient_id and ages of all patients.

# Q2. Select all female patients who are older than 40.

In our patient population, there are no female individuals exceeding the age of 40.



```sql
1 • SELECT * FROM diabetes_prediction
2   where gender = 'Female' and age > 40
```
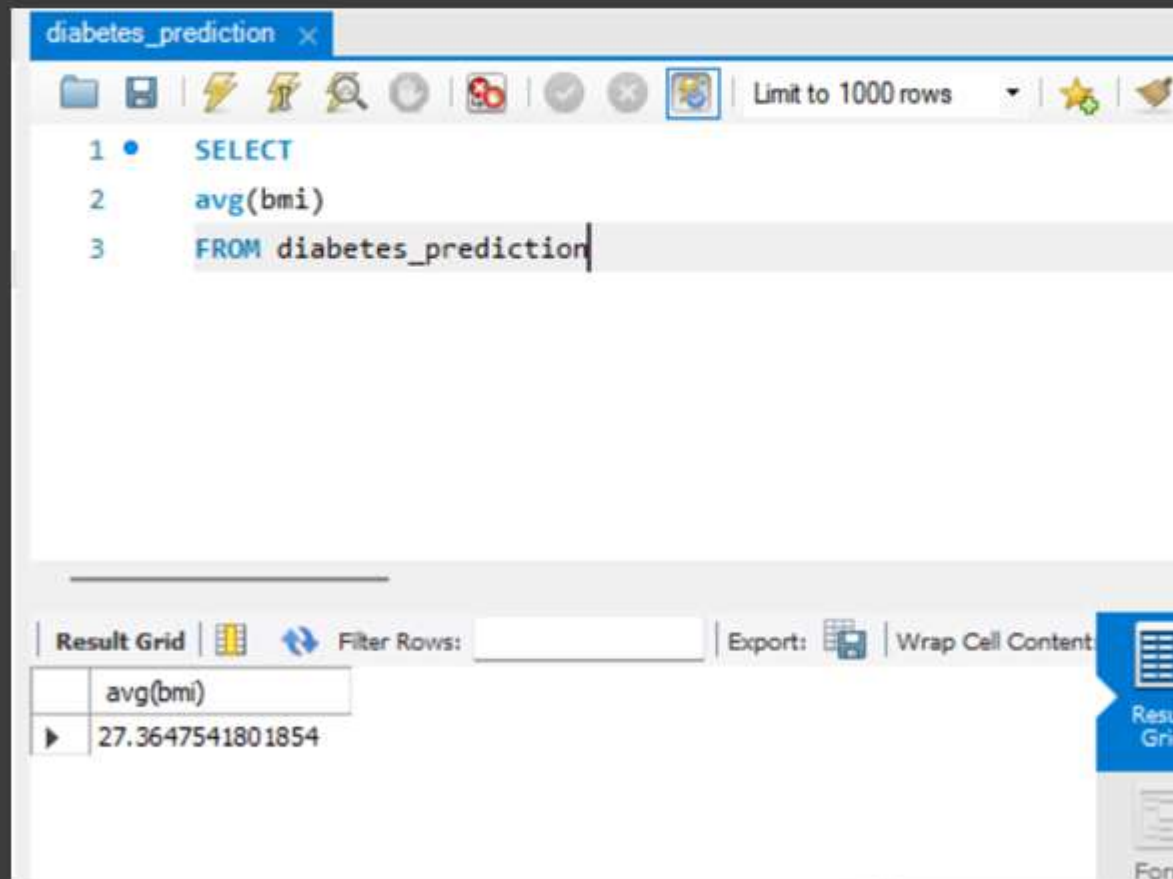
Result Grid | Filter Rows: | Export: | Wrap Cell Content: 

| EmployeeName | Patient_id | gender | D.O.B | hypertension | heart_disease | smoking_history | bmi | HbA1c_level | blood_glu |
|---|---|---|---|---|---|---|---|---|---|

# Q3. Calculate the average BMI of patients.

# Q4. List patients in descending order of blood glucose levels.

# Q5. Find patients who have hypertension and diabetes.

# Q6. Determine the number of patients with heart disease.



```sql
SELECT
count(heart_disease) as number_patients_with_heart_disease
FROM diabetes_prediction
where heart_disease = 1
```

| number_patients_with_heart_disease |
| --- |
| 155 |

# Q7. Group patients by smoking history and count how many smokers and nonsmokers there are.

# Q8. Retrieve the Patient_ids of patients who have a BMI greater than the average BMI.

# Q9. Find the patient with the highest HbA1c level and the patient with the lowest HbA1clevel.

## ⊙ Highest



## ⊙ Lowest

## Q10. Calculate the age of patients in years (assuming the current date as of now).

# Q11. Rank patients by blood glucose level within each gender group.



```sql
select
* , rank() over(partition by gender order by blood_glucose_level) as ranking
from diabetes_prediction
```

| EmployeeName | Patient_id | gender | D.O.B | hypertension | heart_disease | smoking_history | bmi | HbA1c_level | blood_glucose_level | diabetes | age | ranking |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CARLOS RECINOS | PT2286 | Female | 4/16/1999 | 0 | 0 | No Info | 16.3 | 4.8 | 80 | 0 | 24 | 1 |
| SANDY CUADRA | PT3938 | Female | 5/30/1999 | 0 | 0 | No Info | 26.26 | 6 | 80 | 0 | 24 | 1 |
| DERMOT DORGAN | PT2420 | Female | 4/20/1999 | 0 | 0 | No Info | 27.32 | 5.7 | 80 | 0 | 24 | 1 |
| JOSEPH ROBLES | PT3909 | Female | 5/29/1999 | 0 | 0 | never | 22.61 | 6.6 | 80 | 0 | 24 | 1 |
| SARAH WILNER | PT2402 | Female | 4/19/1999 | 0 | 0 | No Info | 16.7 | 5.7 | 80 | 0 | 24 | 1 |
| LORI CADIGAN | PT2235 | Female | 4/15/1999 | 1 | 0 | never | 28.62 | 6.1 | 80 | 0 | 24 | 1 |
| SANJAI NATH | PT2952 | Female | 5/3/1999 | 0 | 0 | never | 18.38 | 5.8 | 80 | 0 | 24 | 1 |
| JERRY TIDWELL | PT3697 | Female | 5/24/1999 | 0 | 0 | No Info | 22.58 | 6 | 80 | 0 | 24 | 1 |
| MICHELLE JEAN | PT2365 | Female | 4/18/1999 | 0 | 0 | ever | 35.83 | 6.6 | 80 | 0 | 24 | 1 |
| STEPHANIE STUART | PT3432 | Female | 5/18/1999 | 0 | 0 | former | 27.74 | 5.7 | 80 | 0 | 24 | 1 |
| TAIRA DE BERNARDI | PT3672 | Female | 5/23/1999 | 0 | 0 | No Info | 27.32 | 6.6 | 80 | 0 | 24 | 1 |

# Q12. Update the smoking history of patients who are older than 50 to "Ex-smoker."

# Q13. Insert a new patient into the database with sample data.



```sql
Insert INTO diabetes_prediction
(EmployeeName, Patient_id, gender, age, hypertension, heart_disease, smoking_history, bmi, HbA1c_level, blood_glucose_level, diabetes)
values("Ram Sinha", 'PT1000000', "Male", "02/14/2022", 0, 0, "never", 25.19, 5,80, 0, 22);
select * from diabetes_prediction
where Patient_id= 'PT1000000';
```
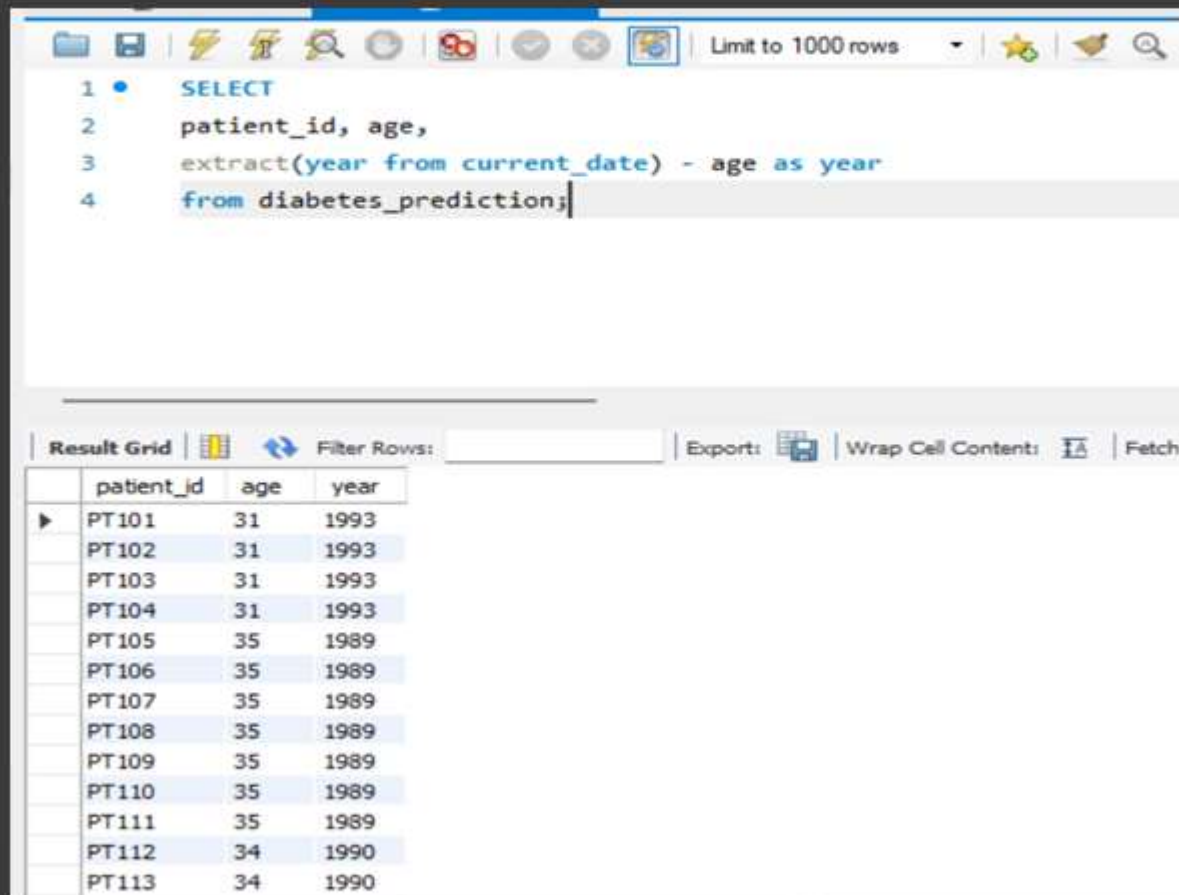
| EmployeeName | Patient_id | gender | D.O.B | hypertension | heart_disease | smoking_history | bmi | HbA1c_level | blood_glucose_level | diabetes | age |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Ram Sinha | PT1000000 | Male | 02/14/2022 | 0 | 0 | never | 25.19 | 5 | 80 | 0 | 22 |

# Q14. Delete all patients with heart disease from the database.



```
diabetes_prediction  ×

1 ●    delete from  diabetes_prediction
2      where heart_disease = 1;
3 ●    select * from diabetes_prediction
```

# Q15. Find patients who have hypertension but not diabetes using the EXCEPT operator.



```
1 •    select * from diabetes_prediction
2      where hypertension = 1
3 ☒    Except
4      select * from diabetes_prediction
5      where diabetes = 1
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content:

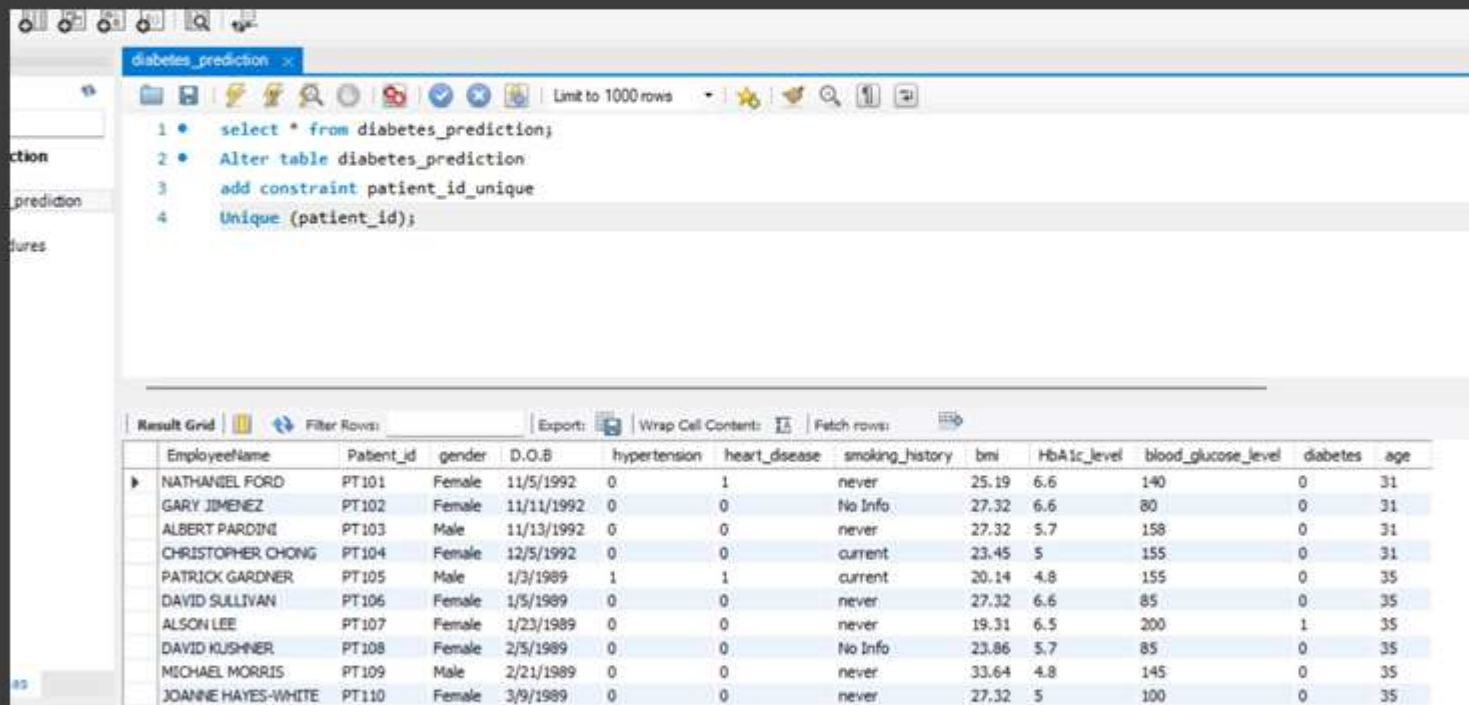| EmployeeName | Patient_id | gender | D.O.B | hypertension | heart_disease | smoking_history | bmi | HbA1c_level | blood_glucose_level | diabetes | age |
|---|---|---|---|---|---|---|---|---|---|---|---|
| PATRICK GARDNER | PT105 | Male | 1/3/1989 | 1 | 1 | current | 20.14 | 4.8 | 155 | 0 | 35 |
| DENISE SCHMITT | PT129 | Male | 6/29/1989 | 1 | 0 | never | 26.47 | 4 | 158 | 0 | 34 |
| THOMAS SIRAGUSA | PT143 | Female | 9/16/1989 | 1 | 1 | never | 32.02 | 5 | 159 | 0 | 34 |
| RAY CRAWFORD | PT155 | Female | 1/2/1997 | 1 | 0 | never | 23.05 | 4.8 | 130 | 0 | 27 |
| KENNETH SMITH | PT161 | Male | 3/9/1997 | 1 | 0 | current | 27.86 | 6.6 | 145 | 0 | 27 |
| CHARLES SCOTT | PT215 | Female | 6/8/1997 | 1 | 0 | never | 34.2 | 5.7 | 140 | 0 | 26 |
| LESLIE DUBBIN | PT220 | Male | 6/19/1997 | 1 | 1 | current | 27.32 | 5 | 126 | 0 | 26 |
| SHANNON SAKOWSKI | PT227 | Male | 7/2/1997 | 1 | 0 | No Info | 28.73 | 6.6 | 160 | 0 | 26 |
| MARISA MORET | PT241 | Female | 7/13/1997 | 1 | 0 | never | 44.06 | 6.5 | 160 | 0 | 26 |
| STEPHEN TACCHINI | PT326 | Female | 8/28/1997 | 1 | 0 | never | 36.73 | 6.6 | 126 | 0 | 26 |
| ANDREW LOGAN | PT339 | Male | 9/5/1997 | 1 | 0 | No Info | 25.31 | 6 | 130 | 0 | 26 |

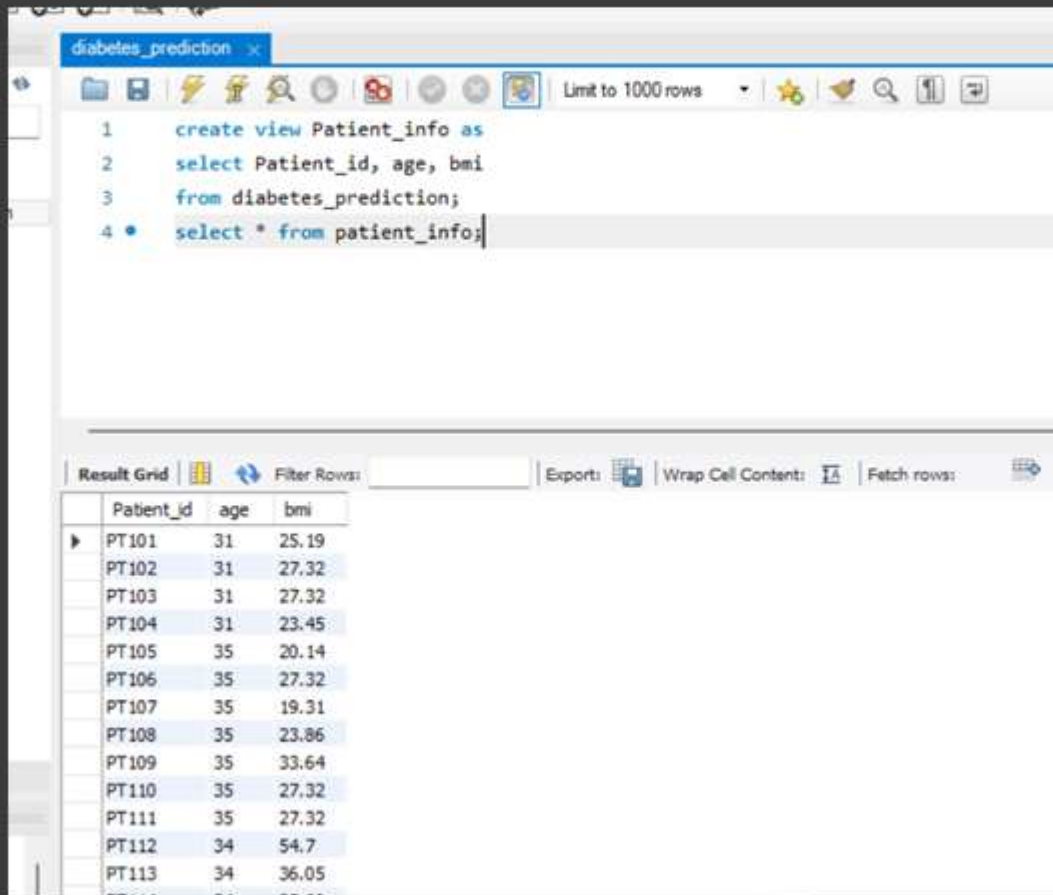# Q16. Define a unique constraint on the "patient_id" column to ensure its values are unique.

# Q17. Create a view that displays the Patient_ids, ages, and BMI of patients.

## Q18. Suggest improvements in the database schema to reduce data redundancy and improve data integrity.

1. **Normalization (Up to 3NF):** Organize data into separate tables based on functional dependencies to eliminate redundancy and ensure data integrity.

2. **Primary Keys and Foreign Keys:** Use primary keys to uniquely identify rows in each table and foreign keys to establish relationships between tables, enforcing referential integrity.

3. **Avoid Multi-Valued Dependencies:** Identify and eliminate multi-valued dependencies by decomposing tables or creating separate tables for related attributes to maintain atomicity and data integrity.

4. **Denormalization for Performance**: Consider denormalizing tables or introducing calculated fields to improve query performance, but carefully balance performance gains with the risk of data inconsistency.

# Q19. Explain how you can optimize the performance of SQL queries on this dataset.

1. **Indexing:** Use indexes on columns frequently used in WHERE clauses, JOIN conditions, and ORDER BY clauses to expedite row retrieval.

2. **Optimize Joins:** Prefer INNER JOIN over OUTER JOIN for efficiency. Prioritize smaller tables as driving tables in the join clause.

3. **Avoid SELECT :** Explicitly specify needed columns instead of using SELECT *, reducing data retrieval overhead.

4. **Limit Result Set:** Apply LIMIT clause to restrict the number of rows returned, minimizing data transfer and processing.