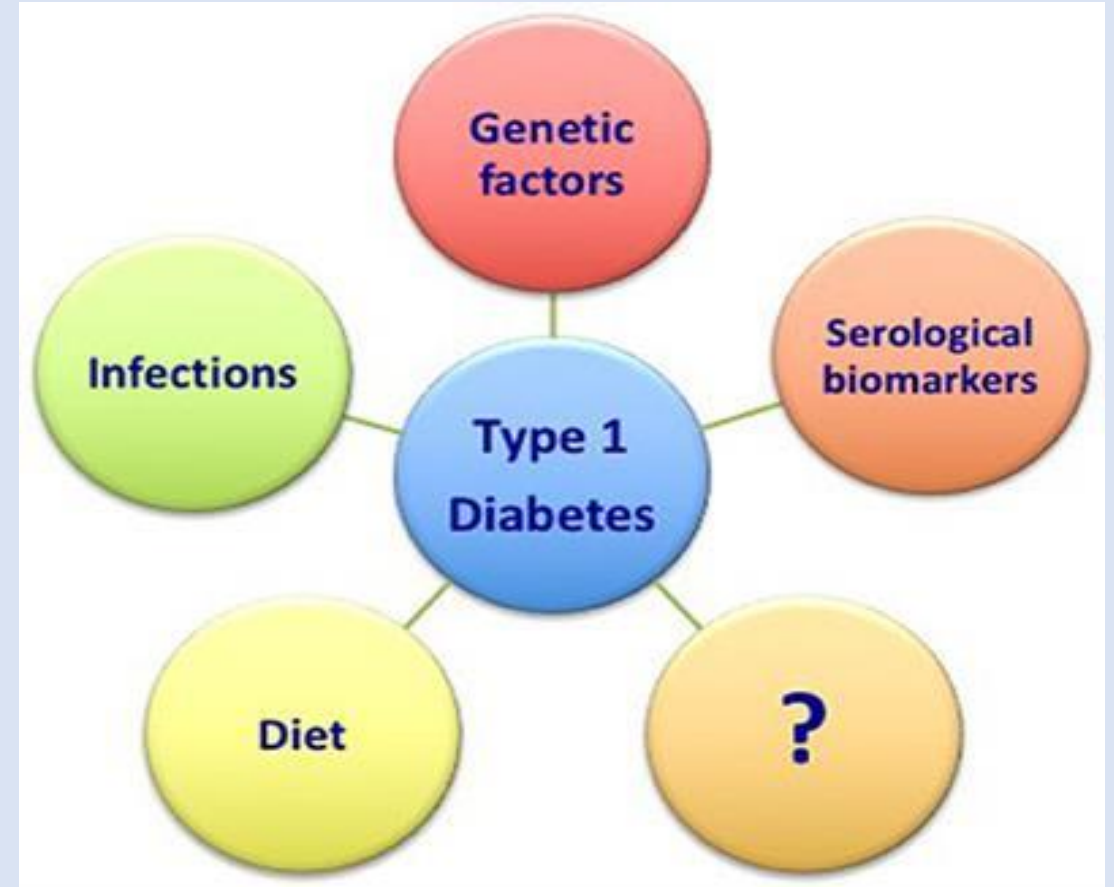


# DIABETES PREDICTION ANALYSIS



# Data Acquisition

- This Analysis was conducted using a comprehensive dataset provided by **PSYLIQ**.

## Project Highlights

- A detailed analysis was conducted for Patients with diabetes based on the following factors:
- Age, smoking history, heart disease and so on .

## Project Tool

- SQL (BigQuery)

# Data Cleaning Process

- Checked for Duplicates, found None.

```
1  -- Checking for Duplicates
2  SELECT
3  | EmployeeName,Patient_id,gender,age,hypertension,heart_disease,smoking_history,bmi,HbA1c_level, blood_glucose_level, diabetes
4  FROM `psyliq-intern.Diabetes_prediction.Diabetes 2`
5  GROUP BY
6  | EmployeeName,Patient_id,gender,age, hypertension,heart_disease,smoking_history,bmi,HbA1c_level, blood_glucose_level,diabetes
7  HAVING COUNT(*) > 1;
8
9
```

Query results

JOB INFORMATION

RESULTS

CHART

PREVIEW

JSON

EXECUTION DETAILS

EXECUTION GRAPH



There is no data to display.

Checked for Null Values, found none

```
1  -- Checking for Null Values
2  SELECT
3  | EmployeeName, Patient_id, gender, age, hypertension, heart_disease, smoking_history, bmi, HbA1c_level, blood_glucose_level, diabetes
4  FROM `psyliq-intern.Diabetes_prediction.Diabetes` 2
5  WHERE EmployeeName IS NULL
6  OR Patient_id IS NULL
7  OR gender IS NULL
8  OR age IS NULL
9  OR Hypertension IS NULL
10 OR heart_disease IS NULL
11 OR smoking_history IS NULL
12 OR bmi IS NULL
13 OR HbA1c_level IS NULL
14 OR blood_glucose_level IS NULL
15 OR diabetes IS NULL;
16
17
18
```

### Query results

JOB INFORMATION

RESULTS

CHART

PREVIEW

JSON

EXECUTION DETAILS

EXECUTION GRAPH



There is no data to display.

1). Retrieve the Patient\_id and ages of all patients.

```
1  -- Retrieve Patient_id and Ages of all patients
2  SELECT
3      EmployeeName,
4      Patient_id,
5      CAST(age AS INT64) AS Age
6  FROM `psyliq-intern.Diabetes_prediction.Diabetes_2`
7  ORDER BY
8      Patient_id,
9      Age
```

### Query results

JOB INFORMATION		RESULTS	CHART	PREVIEW	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	EmployeeName	Patient_id	Age				
1	SHERYL BREGMAN	PT1000	34				
2	JOHN HOFFMAN	PT10000	46				
3	Meredith H Reddoch-Ho	PT100000	45				
4	Minouche Kandel	PT100001	57				
5	Mose Thornton	PT100002	72				
6	Helen H Chong	PT100003	44				
7	Marvin M Mouton	PT100004	62				
8	Edward A Ang	PT100005	21				
9	Gordon G Leong	PT100006	62				
10	Judith Reyes	PT100007	3				
11	Tara L Croan	PT100008	48				
12	Brian M DeNave	PT100009	67				
13	ROBERTO VALLADARES	PT10001	27				
14	Jennifer J Pascual	PT100010	41				

2. Select all female patients who are older than 40.

```
1  -- Female Patients older than 40
2  SELECT
3      EmployeeName,
4      Gender,
5      CAST(age AS INT64) AS Age
6  FROM `psyliq-intern.Diabetes_prediction.Diabetes_2`
7  WHERE
8      Gender = 'Female'
9      AND Age >40
10
```

#### Query results

JOB INFORMATION		RESULTS	CHART	PREVIEW	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	EmployeeName	Gender	Age				
1	FRANCIS HAGAN	Female	42				
2	MARIALUZ BANARES	Female	42				
3	RICHARD ACERET	Female	42				
4	MICHAEL MITCHELL	Female	42				
5	AMOD DHAKAL	Female	42				
6	RENE QUIAMBAO	Female	42				
7	CHARLES SHEEHAN	Female	42				
8	DAMIAN DAVIS	Female	42				
9	ANNE KOFMAN	Female	42				
10	CORY DECKER	Female	42				
11	DANIEL IP	Female	42				
12	DANIELLE THOMAS	Female	42				
13	VIRGILIO RAFANAN	Female	42				
14	YVETTE WILLIAMS-DUBRIWNY	Female	42				

3). Calculate the average BMI of patients.

```
1  -- Average BMI of Patients
2  SELECT
3  |  CAST(AVG(bmi) AS INT64) AS Average_BMI
4  FROM `psyliq-intern.Diabetes_prediction.Diabetes 2`
```

### Query results

JOB INFORMATION

RESULTS

CHART

PREVIEW

JSON

EXECUTION DETAILS

EXECUTION GRAPH

Row	Average_BMI	
1	27	

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

```
1  -- Blood glucose level of patients in descending order
2  SELECT
3      EmployeeName,
4      Patient_id,
5      Blood_glucose_level
6  FROM `psyliq-intern.Diabetes_prediction.Diabetes_2`
7  ORDER BY
8      Blood_glucose_level
9      DESC
10
11
```

#### Query results

JOB INFORMATION		RESULTS	CHART	PREVIEW	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	EmployeeName	Patient_id	Blood_glucose_level				
1	Carlos Imbong	PT51740	300				
2	Santiago Lagandaon	PT49657	300				
3	Yu Qiong Wen	PT63081	300				
4	Rita L Kearns	PT76623	300				
5	STEPHEN ZOLLMAN	PT5536	300				
6	Rosa Narvaez	PT57767	300				
7	Teresa Chan	PT45322	300				
8	Viet Q Ha	PT78543	300				
9	KATHLEEN GULBENGAY	PT36031	300				
10	Michael Callejas	PT59205	300				
11	Jeff Fisher	PT54591	300				
12	Theresa Smith	PT65140	300				
13	Sergey Trofimenko	PT89757	300				
14	Flor D Roman	PT99809	300				



5). Find patients who have hypertension and diabetes.

```
1  -- Patients who have hypertension and diabetes
2  SELECT
3      EmployeeName,
4      Hypertension,
5      Diabetes
6  FROM `psyliq-intern.Diabetes_prediction.Diabetes` 2
7  WHERE
8      Hypertension <> 0
9      AND Diabetes <> 0
10
```

#### Query results

JOB INFORMATION		RESULTS	CHART	PREVIEW	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	EmployeeName	Hypertension	Diabetes				
1	Richard S Burns	1	1				
2	Matelina Alexander	1	1				
3	Jeffrey Black	1	1				
4	Anthony Branchcomb	1	1				
5	Jason Jefferson	1	1				
6	Irving McKnight	1	1				
7	Joseph B Rocero	1	1				
8	Richard Sucre	1	1				
9	Justin Morgan	1	1				
10	ALISON YEE	1	1				
11	Xiao Dong Cai	1	1				
12	Jeffrey E Airth	1	1				
13	Barbara L Moss	1	1				
14	Randy F Acosta	1	1				
15	MICHELLE SHINN	1	1				

6). Determine the number of patients with heart disease.

➔ The number of patients with heart disease in this dataset is 3,942.

```
1  -- Number of Patients with Heart Disease
2  SELECT
3  | COUNT(Heart_Disease) AS Heart_Disease
4  FROM `psyliq-intern.Diabetes_prediction.Diabetes 2`
5  WHERE
6  | Heart_Disease <> 0
7
```

### Query results

JOB INFORMATION

RESULTS

CHART

PREVIEW

JSON

EXECUTION DETAILS

EXECUTION GRAPH

Row	Heart_Disease
1	3942

7). Group patients by smoking history and count how many smokers and non-smokers there are.

```
1 --Group patients by smoking history and count how many smokers and non-smokers there are.
2 select Smoking_history,COUNT(Smoking_history) As Total_Smoking_history
3 | from 'psyliq-intern.Diabetes_prediction.Diabetes 2'
4 | where Smoking_history is not null
5 group by Smoking_history
```

### Query results

JOB INFORMATION		RESULTS	CHART	PREVIEW	JSON	EXECUTION DETAILS	E
Row	Smoking_history	Total_Smoking_history					
1	never	35095					
2	No Info	35816					
3	not current	6447					
4	current	9286					
5	ever	4004					
6	former	9352					

8) Retrieve the Patient\_ids of patients who have a BMI greater than the average BMI.

```
1  -- Patients with greater than average BMIs
2  -- Average BMI is 27
3  SELECT
4    Patient_id,
5    CAST(bmi AS int64) AS BMI
6  FROM `psyliq-intern.Diabetes_prediction.Diabetes` 2
7  Where
8    BMI > 27
```

#### Query results

JOB INFORMATION		RESULTS	CHART	PREVIEW	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	Patient_id	BMI					
1	PT50982	57					
2	PT4330	57					
3	PT22555	88					
4	PT41474	49					
5	PT13977	51					
6	PT68256	48					
7	PT97017	51					
8	PT59611	52					
9	PT90093	49					
10	PT96442	51					
11	PT72360	48					
12	PT5314	50					
13	PT40882	51					

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

➔ Patients with HbA1c level 9 were the highest and Patients with HbA1c level 3.5 were the lowest as shown below:

```
1 -- Patients with the Highest HbA1c level
2 SELECT Patient_id, MAX(HbA1c_level) AS Highest_HbA1c_level
3 from `psyliq-intern.Diabetes_prediction.Diabetes_2`
4 where patient_id is not null
5 Group by Patient_id
6 Order by Highest_HbA1c_level DESC
7
8
9
10
```

#### Query results

JOB INFORMATION		RESULTS	CHART	PREVIEW	JSON
Row	Patient_id	Highest_HbA1c_level			
1	PT40453	9.0			
2	PT44779	9.0			
3	PT53598	9.0			
4	PT71476	9.0			
5	PT73983	9.0			
6	PT36071	9.0			
7	PT59549	9.0			
8	PT11499	9.0			
9	PT11719	9.0			
10	PT1866	9.0			
11	PT87900	9.0			
12	PT56655	9.0			
13	PT78316	9.0			

```
1 -- Patients with the Lowest HbA1c level
2 SELECT Patient_id, MIN(HbA1c_level) AS Lowest_HbA1c_level
3 from `psyliq-intern.Diabetes_prediction.Diabetes_2`
4 where patient_id is not null
5 Group by Patient_id
6 Order by Lowest_HbA1c_level ASC
7
8
9
10
```

#### Query results

JOB INFORMATION		RESULTS	CHART	PREVIEW	JSON
Row	Patient_id	Lowest_HbA1c_level			
1	PT891	3.5			
2	PT1141	3.5			
3	PT1841	3.5			
4	PT1940	3.5			
5	PT2032	3.5			
6	PT5457	3.5			
7	PT6806	3.5			
8	PT7552	3.5			
9	PT9032	3.5			
10	PT9034	3.5			
11	PT10302	3.5			
12	PT10389	3.5			
13	PT11114	3.5			

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

```
1 -- Calculate the age of patients in years (assuming the current date as of now).
2 SELECT Patient_id,
3 Cast (Age_ AS int64) AS Age
4 FROM 'psyliq-intern.Diabetes_prediction.Diabetes_2'
5 Order by Age DESC
```

#### Query results

JOB INFORMATION		RESULTS	CHART	PREVIEW	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	Patient_id	Age					
2	PT91186	80					
3	PT65525	80					
4	PT22551	80					
5	PT59430	80					
6	PT4009	80					
7	PT85234	80					
8	PT39672	80					
9	PT60419	80					
10	PT36000	80					
11	PT60494	80					
12	PT78666	80					
13	PT1066	80					
14	PT76319	80					
15	PT83271	80					
16	PT86928	80					
17	PT72902	80					
18	PT92271	80					
19	PT94953	80					

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

```
1  --Rank patients by blood glucose level within each gender group
2  with blood_glucose_level as(
3      select Gender, Blood_glucose_level, Patient_id
4      from `psyliq-intern.Diabetes_prediction.Diabetes_2`
5      group by Blood_glucose_level, Gender, Patient_id
6  )
7  SELECT
8      Patient_id, Blood_glucose_level, Gender, ROW_NUMBER() Over(Order by Blood_glucose_level DESC) as Rank
9  FROM `psyliq-intern.Diabetes_prediction.Diabetes_2`
10
```

#### Query results

JOB INFORMATION							RESULTS	CHART	PREVIEW	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	Patient_id	Blood_glucose_level	Gender	Rank								
1	PT51740	300	Female	1								
2	PT49657	300	Male	2								
3	PT63081	300	Female	3								
4	PT76623	300	Female	4								
5	PT55536	300	Female	5								
6	PT57767	300	Female	6								
7	PT45322	300	Female	7								
8	PT78543	300	Female	8								
9	PT36031	300	Female	9								
10	PT59205	300	Male	10								
11	PT54591	300	Female	11								
12	PT65140	300	Female	12								
13	PT89757	300	Male	13								
14	PT99809	300	Male	14								

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

```
1  --Update the smoking history of patients who are older than 50 to "Ex-smoker."  
2  update 'psyliq-intern.Diabetes_prediction.Diabetes 2'  
3  Set smoking_history = "Ex-smoker"  
4  where age > 5.0  
5
```



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

```
1  --Insert a new patient into the database with sample data.  
2  insert into 'psyliq-intern.Diabetes_prediction.Diabetes 2'  
3  values("Jordin", PT5654, Female, 5.0, 0, 0, "Ex-smoker", 22.65, 9, 155, 0)
```

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

```
1  --Delete all patients with heart disease from the database
2  DELETE from`psyliq-intern.Diabetes_prediction.Diabetes` 2`
3  where Heart_disease <> 0
4
```

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

```
1  -- Find patients who have hypertension but not diabetes using the EXCEPT operator.
2  SELECT patient_id, hypertension, diabetes
3  FROM psyliq-intern.Diabetes_prediction.Diabetes 2
4  WHERE hypertension = 1 AND diabetes = 0
```

### Query results

JOB INFORMATION		RESULTS	CHART	PREVIEW	JSON	EXECUTION DETAILS	EXECUTION GRAPH
Row	patient_id	hypertension	diabetes				
2	PT87841	1	0				
3	PT92208	1	0				
4	PT59611	1	0				
5	PT994	1	0				
6	PT14638	1	0				
7	PT5413	1	0				
8	PT4435	1	0				
9	PT74700	1	0				
10	PT44273	1	0				
11	PT79212	1	0				
12	PT24713	1	0				
13	PT18824	1	0				
14	PT33498	1	0				
15	PT80637	1	0				
16	PT76198	1	0				
17	PT72205	1	0				
18	PT1186	1	0				
19	PT18676	1	0				

16). Define a unique constraint on the "patient\_id" column to ensure its values are unique.

➔ A unique constraint on the Patient\_id column would be to ensure its not Null to be sure its values are unique.

1.17). Create a view that displays the Patient\_ids, ages, and BMI of patients.

➔ After running a query, click the **Save view** button above the query results window to save the query as a view.

In the **Save view** dialog:

For **Project name**, select a project to store the view.

For **Dataset name**, choose a dataset to store the view. The dataset that contains your view and the dataset that contains the tables referenced by the view must be in the same [location](#).

For **Table name**, enter the name of the view.

Click **Save**.

Save view

❗ The destination dataset for a saved view must be in the same region as the source, otherwise a 'Dataset not found' error will be returned.

Project \*  
psyliq-intern BROWSE

Dataset \*  
Diabetes\_prediction

Table \*  
Diabetes\_2

Maximum name size is 1,024 UTF-8 bytes. Unicode letters, marks, numbers, connectors, dashes and spaces are allowed. The job will create the specified destination table if needed.

SAVE CANCEL

SCHEMA

DETAILS

LINEAGE

DATA PROFILE

DATA QUALITY

Filter

<input type="checkbox"/>	Field name	Type	Mode	Key	Collation	Default value	Policy tags <sup>1</sup>	Description
<input type="checkbox"/>	patient_id	STRING	NULLABLE	-	-	-	-	-
<input type="checkbox"/>	Age	INTEGER	NULLABLE	-	-	-	-	-
<input type="checkbox"/>	BMI	INTEGER	NULLABLE	-	-	-	-	-

EDIT SCHEMA

SCHEMA DETAILS LINEAGE DATA PROFILE DATA QUALITY

### View info

EDIT DETAILS

View ID	psyliq-intern.Diabetes_prediction.Diabetes_2
Created	9 Feb 2024, 10:00:42 UTC+1
Last modified	9 Feb 2024, 10:00:48 UTC+1
View expiry	9 Apr 2024, 10:00:42 UTC+1
Use Legacy SQL	false
Policy	
Privacy unit column	
Threshold	
Description	
Labels	
Primary key(s)	
Tags	

### Storage info

Number of rows	0
Total logical bytes	0 B
Active logical bytes	0 B
Long-term logical bytes	0 B
Total physical bytes	0 B
Active physical bytes	0 B
Long-term physical bytes	0 B
Time travel physical bytes	0 B

### Query

EDIT QUERY

```
1 -- Create a view that displays the Patient_ids, ages, and BMI of patients.
2 SELECT patient_id, Age_, Cast(BMI As int64) AS BMI
3 FROM `psyliq-intern.Diabetes_prediction.Diabetes_2`
4
```

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

➔ The suggestions below can help improve the database schema to reduce data redundancy and improve data integrity:

- Data integrity mechanisms are essential for reducing data redundancy. By implementing constraints and triggers, you can prevent duplicate data from being entered, enforce referential integrity, and streamline data retrieval thereby ensuring data accuracy, consistency, and validity.
- Also, regular reviewing and updating data integrity rules as your data evolves, helps to maintain data quality and prevent inconsistencies.

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

To optimize the performance of SQL queries on this dataset, the following can be done:

- Minimize the use of subqueries.
- Avoid unnecessary data retrieval.
- Retrieve only necessary columns.
- SELECT fields instead of using SELECT \*.

Thank You