**Project Title:** Smoking and Cancer Risk Analysis using MY SQL

**Name:** Venkateswarlu Pujari

**Date:** 18th-04- 2025

**Summary:**

In this project, MySQL was used to analyse a healthcare dataset focused on understanding how smoking habits correlate with cancer risk. The data was organized across three relational tables—Patients, Habits, and Results—covering lifestyle, smoking behaviour, and medical outcomes. Through a series of SQL queries, the analysis identified high-risk groups, examined smoking patterns, calculated BMI trends, and evaluated cancer distribution. Key insights included patterns in cancer types by gender, risk levels based on smoking history, and cases of long-term smokers without diagnosis—offering valuable implications for preventive healthcare and research.

**Query Analysis**

1.  **Query:** List all patients who are older than 50 and have a poor diet.

```
10
11 ●    SELECT
12          *
13      FROM
14          project1.patients
15      WHERE
16          age > 50 AND diet_quality = 'Poor';
```

**Output:**

| Patient_ID | Age | Gender | BMI | Physical_Activity_Level | Diet_Quality |
|---|---|---|---|---|---|
| 17 | 75 | Other | 21.4 | Low | Poor |
| 21 | 77 | Male | 28.4 | Moderate | Poor |
| 26 | 68 | Male | 27 | Low | Poor |
| 27 | 81 | Male | 35.6 | High | Poor |
| 31 | 56 | Female | 29.7 | Low | Poor |
| 34 | 77 | Female | 20 | Moderate | Poor |
| 39 | 88 | Female | 21.5 | Moderate | Poor |
| 46 | 71 | Female | 23.4 | Low | Poor |
| 49 | 80 | Female | 25.8 | Moderate | Poor |
| 55 | 65 | Other | 21.7 | Moderate | Poor |
| 56 | 89 | Male | 17.7 | High | Poor |
| 66 | 88 | Other | 19.2 | High | Poor |
| 78 | 89 | Male | 30.1 | High | Poor |

patients 1 ✕

**Insight**: These are high-risk patients based on age and dietary habits.

2.  **Query:** Display the distinct physical activity levels recorded in the dataset

```
12
13 ●    SELECT DISTINCT
14          (physical_activity_level)
15      FROM
16          project1.patients;
17
```

**Output:**

| physical_activity_level |
|---|
| ▸ Low |
| High |
| Moderate |

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

**Insight:** The dataset includes distinct Low, Moderate, and High activity levels.

3. **Query:** Show the top 5 patients with the highest number of cigarettes smoked per day

```
14
15 •   SELECT
16         *
17     FROM
18         project1.habits
19     ORDER BY cigarettes_per_day DESC
20     LIMIT 5;
21
```

**Output:**

Result Grid | Filter Rows: | Export: | Wrap Cell Content: 𝐼𝐴 | Fetch rows:

| Patient_ID | Smoking_Status | Cigarettes_Per_Day | Years_Smoking |
|---|---|---|---|
| ▸ 590 | Former | 20 | 39 |
| 1730 | Former | 20 | 2 |
| 1396 | Former | 20 | 16 |
| 1627 | Former | 20 | 16 |
| 1794 | Former | 19 | 15 |

**Insight:** These patients have the highest daily tobacco exposure.

4. **Query:** Find the average BMI grouped by physical activity level.

```
20
21 •   SELECT
22         physical_activity_level, ROUND(AVG(bmi), 2) AS Avg_BMI
23     FROM
24         project1.patients
25     GROUP BY physical_activity_level;
26
```

**Output:**

| physical_activity_level | Avg_BMI |
|---|---|
| Low | 25.1 |
| High | 24.97 |
| Moderate | 25 |

**Insight:** Patients with high physical activity levels tend to have lower BMI on average.

5. **Query:** Count how many patients fall into each smoking status category.

```
22
23 •    SELECT
24          COUNT(patient_id) AS Total_Patients, smoking_status
25      FROM
26          project1.habits
27      GROUP BY smoking_status;
28
```

**Output:**

| Total_Patients | smoking_status |
|---|---|
| 1054 | Never |
| 779 | Former |
| 734 | Current |

**Insight:** Most patients have never smoked. The remaining patients are more former smokers than current smokers.

6. **Query:** List the number of cancer patients per cancer type ordered by patient count

```
24
25 •    SELECT
26          cancer_type, COUNT(patient_id) AS Total_Patients
27      FROM
28          project1.results
29      GROUP BY cancer_type
30      ORDER BY Total_Patients DESC;
31
```

**Output:**

**Insight:** The highest concentration of patients falls under the "None" category.

7. **Query:** Find the average years of smoking for patients who are current and former smokers.

```
26
27 •  SELECT
28         smoking_status, AVG(years_smoking) AS Avg_years_of_smoking
29     FROM
30         project1.habits
31     GROUP BY smoking_status
32     HAVING smoking_status IN ('Current' , 'former');
33
```

**Output:**



| smoking_status | Avg_years_of_smoking |
|---|---|
| Former | 25.1579 |
| Current | 24.9360 |

**Insight:** Former smokers have a slightly longer smoking history compared to current smokers

8. **Query:** Get the gender-wise average BMI of patients with cancer.

```
29
30 •    SELECT
31         patients.gender, ROUND(AVG(patients.bmi), 2) AS Avg_BMI
32     FROM
33         project1.patients
34             INNER JOIN
35         project1.results ON patients.patient_id = results.patient_id
36     WHERE
37         results.cancer_type != 'none '
38     GROUP BY patients.gender;
39
```

**Output:**

| | gender | Avg_BMI |
|---|---|---|
| ▶ | Other | 24.87 |
| | Male | 25.13 |
| | Female | 25.08 |

**Insight:** Males have the highest average BMI compared to females and other gender categories

9. **Query:** List patients who smoke more than the average number of cigarettes per day.

```
40
41 •  SELECT
42         patient_id,
43         ROUND(AVG(cigarettes_per_day), 0) AS Avg_No_Of_Cigerettes_Per_Day
44     FROM
45         project1.habits
46     WHERE
47         smoking_status != 'Never'
48     GROUP BY patient_id
49   ⊖ HAVING AVG(cigarettes_per_day) > (SELECT
50                 AVG(cigarettes_per_day)
51         FROM
52             project1.habits
53         WHERE
54             smoking_status != 'never');
55
```

**Output:**

| | patient_id | Avg_No_Of_Cigerettes_Per_Day |
|---|---|---|
| ▶ | 3 | 11 |
| | 11 | 15 |
| | 13 | 12 |
| | 14 | 14 |
| | 16 | 14 |
| | 17 | 13 |
| | 18 | 16 |
| | 20 | 11 |

**Insight:** This outcome shows the average number of cigarettes smoked per day by Patients

10. **Query:** Create a new column named Risk_Level using CASE:
   - "High" if smoking status is 'Current'
   - "Medium" if smoking status is 'Former'
   - Else "Low"

```
43
44 •   select *,
45 ⊖   case
46       when smoking_status = "Never" then "Low"
47       when smoking_status = "current" then "High"
48       else "Medium"
49       end as Risk_Level from project1.habits;
50
```

**Output:**

| Patient_ID | Smoking_Status | Cigarettes_Per_Day | Years_Smoking | Risk_Level |
|---|---|---|---|---|
| 1 | Never | 0 | 0 | Low |
| 2 | Never | 0 | 0 | Low |
| 3 | Former | 11 | 44 | Medium |
| 4 | Never | 0 | 0 | Low |
| 5 | Never | 0 | 0 | Low |
| 6 | Never | 0 | 0 | Low |
| 7 | Never | 0 | 0 | Low |
| 8 | Never | 0 | 0 | Low |
| 9 | Never | 0 | 0 | Low |
| 10 | Never | 0 | 0 | Low |
| 11 | Former | 15 | 13 | Medium |

**Insight:** Patients who have never smoked exhibit a low risk level while former smokers have medium level and current smokers are high risk level

11. **Query:** Identify patients who have smoked for more than 20 years but have not been diagnosed with any cancer.

```
52
53 •   SELECT
54         habits.patient_id,
55         habits.smoking_status,
56         habits.cigarettes_per_day,
57         habits.years_smoking,
58         results.cancer_type
59     FROM
60         project1.habits
61             JOIN
62         project1.results ON habits.patient_id = results.patient_id
63     WHERE
64         years_smoking > 20
65             AND cancer_type = 'none';
66
```

**Output:**

| patient_id | smoking_status | cigarettes_per_day | years_smoking | cancer_type |
|---|---|---|---|---|
| 3 | Former | 11 | 44 | None |
| 13 | Current | 12 | 27 | None |
| 16 | Current | 14 | 37 | None |
| 18 | Current | 16 | 41 | None |
| 21 | Former | 12 | 47 | None |
| 30 | Former | 11 | 24 | None |
| 39 | Current | 6 | 45 | None |
| 41 | Former | 10 | 43 | None |
| 47 | Former | 12 | 21 | None |
| 51 | Current | 9 | 47 | None |

**Insight:** This outcome shows patients who have not been diagnosed with cancer.

**12. Query:** Find the top 3 cancer types with the highest number of female patients.

```
55
56 •    SELECT
57             COUNT(patients.patient_id) AS Total_Patients,
58             patients.gender,
59             results.cancer_type
60        FROM
61             project1.patients
62                 JOIN
63             project1.results ON patients.patient_id = results.patient_id
64        WHERE
65             gender = 'Female'
66        GROUP BY results.cancer_type
67        ORDER BY Total_Patients DESC
68        LIMIT 3;
```

**Output:**

| Total_Patients | gender | cancer_type |
|---|---|---|
| 693 | Female | None |
| 75 | Female | Lung |
| 31 | Female | Other |

**Insight:** The majority of female patients in this dataset have no recorded cancer, while a smaller is diagnosed with lung cancer or other types.