Ashwath J

# AUTO INSURANCE RISK  - BUSINESS REPORT
# GRADED PROJECT

Q1)Write a query to calculate what % of the customers have made a claim in the current exposure period[i.e. in the given dataset]?

**Customer made claim = 34060**

```
3
4   --Q1
5   SELECT count(IDpol)AS TOT_Cust_claimed FROM Auto_insurance
6   WHERE ClaimNb>=1;
7
8   SELECT COUNT(IDpol)AS tot_cust FROM
9   FROM Auto_insurance;
```

| | TOT_Cust_claimed |
|---|---|
| 1 | 34060 |

**Total Customers = 678013**

```
11   SELECT count(IDpol)AS TOT_Cust_claimed FROM Auto_insurance;
```

| | TOT_Cust_claimed |
|---|---|
| 1 | 678013 |

**I.e, 5% of customers have claimed.**

Ashwath J

Q2)
2.1. Create a new column as 'claim_flag' in the table 'auto_insurance_risk' as integer datatype.
2.2. Set the value to 1 when ClaimNb is greater than 0 and set the value to 0 otherwise.

```sql
22    ALTER TABLE Auto_insurance add column claim_flag INT;
23    UPDATE Auto_insurance SET claim_flag = 1
24    WHERE ClaimNb>0;
25    UPDATE Auto_insurance SET claim_flag = 0
26    WHERE ClaimNb<=0;
27
```

| | IDpol | ClaimNb | Exposure | Area | VehPower | VehAge | DrivAge | BonusMalus | VehBrand | VehGas | Density | Region | claim_flag |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 0.1 | D | 5 | 0 | 55 | 50 | B12 | Regular | 1217 | R82 | 1 |
| 2 | 3 | 1 | 0.77 | D | 5 | 0 | 55 | 50 | B12 | Regular | 1217 | R82 | 1 |
| 3 | 5 | 1 | 0.75 | B | 6 | 2 | 52 | 50 | B12 | Diesel | 54 | R22 | 1 |
| 4 | 10 | 1 | 0.09 | B | 7 | 0 | 46 | 50 | B12 | Diesel | 76 | R72 | 1 |
| 5 | 11 | 1 | 0.84 | B | 7 | 0 | 46 | 50 | B12 | Diesel | 76 | R72 | 1 |

Ashwath J

Q3)
3.1. What is the average exposure period for those who have claimed?
3.2. What do you infer from the result? Hint: Use claim_flag variable to group the data.

```
12    --Q3
13    SELECT claim_flag, ROUND (AVG (Exposure), 2) AS AVG_EXP_PERIOD
14    FROM Auto. insurance
15    GROUP BY claim_flag;
16
```

|   | claim_flag | AVG_EXP_PERIOD |
|---|---|---|
| 1 | 0 | 0.52 |
| 2 | 1 | 0.64 |

**Inference:  Avg exposure period is high in case of customers who have claimed.**

Ashwath J

Q4)

4.1. If we create an exposure bucket where buckets are like below, what is the % of total claims by these buckets?

4.2. What do you infer from the summary?

Hint: Buckets are => E1 = 0 to 0.25, E2 = 0.26 to 0.5, E3 = 0.51 to 0.75, E4 > 0.75, You need to consider the ClaimNb field to get the total claim count.

```
18    SELECT bucket_list,Total_Claim_bucketwise,SUM(Total_Claim_bucketwise) OVER()AS Total_Claims
19    FROM(
20    SELECT a.bucket_list,SUM(ClaimNb) AS Total_Claim_bucketwise
21    FROM(
22    SELECT *,
23    CASE
24        WHEN Exposure  between 0 and 0.25 THEN 'E1'
25        WHEN Exposure  between 0.26 and 0.5 THEN 'E2'
26        WHEN Exposure  between 0.51 and 0.75 THEN 'E3'
27        ELSE 'E4'
28    END bucket_list
29    FROM Auto_insurance) a
30    GROUP BY a.bucket_list);
31
```

| | bucket_list | Total_Claim_bucketwise | Total_Claims |
|---|---|---|---|
| 1 | E1 | 7131 | 36102 |
| 2 | E2 | 6481 | 36102 |
| 3 | E3 | 5968 | 36102 |
| 4 | E4 | 16522 | 36102 |

```
17    --Q4
18    SELECT bucket_list,ROUND(CAST(Total_Claim_bucketwise AS FLOAT)/CAST(Total_Claims AS FLOAT)*100,2) AS PERCENT_OF_CLAIMS
19    FROM(
20    SELECT bucket_list,Total_Claim_bucketwise,SUM(Total_Claim_bucketwise) OVER()AS Total_Claims
21    FROM(
22    SELECT a.bucket_list,SUM(ClaimNb) AS Total_Claim_bucketwise
23    FROM(
24    SELECT *,
25    CASE
26        WHEN Exposure  between 0 and 0.25 THEN 'E1'
27        WHEN Exposure  between 0.26 and 0.5 THEN 'E2'
28        WHEN Exposure  between 0.51 and 0.75 THEN 'E3'
```

| | bucket_list | PERCENT_OF_CLAIMS |
|---|---|---|
| 1 | E1 | 19.75 |
| 2 | E2 | 17.95 |
| 3 | E3 | 16.53 |
| 4 | E4 | 45.76 |

- **Highest number of claims is done in exposure of more than 0.75**

Ashwath J

Q5)Which area has the highest number of average claims? Show the data in percentage w.r.t. the number of policies in the corresponding Area.
 Hint: Use the ClaimNb field for this question.

```
36   --Q5
37   SELECT *,ROUND(CAST(Num_of_claims AS Float)/CAST(Num_of_policies AS Float)*100,2) AS percent_of_claim
38   ⊟FROM(
39     SELECT Region,SUM(ClaimNb) AS Num_of_claims,COUNT(IDpol) AS Num_of_policies
40     FROM Auto_insurance
41   └GROUP BY Region)
42     ORDER BY percent_of_claim DESC;
43
```

|   | Region | Num_of_claims | Num_of_policies | percent_of_claim |
|---|--------|---------------|-----------------|------------------|
| 1 | R53    | 2702          | 42122           | 6.41             |
| 2 | R42    | 133           | 2200            | 6.05             |
| 3 | R82    | 5032          | 84752           | 5.94             |
| 4 | R25    | 633           | 10893           | 5.81             |
| 5 | R24    | 9204          | 160601          | 5.73             |

Q6)

If we use these exposure buckets along with Area i.e. group Area and Exposure Buckets together and look at the claim rate, an interesting pattern could be seen in the data. What is that?

```
46      --Q6
47      SELECT bucket_list,Area,ROUND(CAST(num_of_claim AS FLOAT)/CAST(tot_policies AS FLOAT)*100,2) AS claim_rate
48      FROM(
49        SELECT Area,bucket_list,SUM(ClaimNb) AS num_of_claim,COUNT(IDpol) AS tot_policies
50        FROM(
51        SELECT * ,
52        CASE
53            WHEN Exposure  between 0 and 0.25 THEN 'E1'
54            WHEN Exposure  between 0.26 and 0.5 THEN 'E2'
55            WHEN Exposure  between 0.51 and 0.75 THEN 'E3'
56            ELSE 'E4'
57        END bucket_list
58        FROM Auto_insurance)
59        GROUP BY Area,bucket_list)
60        GROUP BY bucket_list,Area
61        ORDER BY Area,bucket_list;
62
```

|   | bucket_list | Area | claim_rate |
|---|-------------|------|-----------|
| 1 | E1 | A | 2.95 |
| 2 | E2 | A | 4.19 |
| 3 | E3 | A | 5.57 |
| 4 | E4 | A | 6.15 |
| 5 | E1 | B | 3.11 |

**INFERENCE- For each area the claim_rate is increasing with the increase in the emposure bucket level.**

Ashwath J

Q7)

7.1. If we look at average Vehicle Age for those who claimed vs those who didn't claim, what do you see in the summary? (1.5+1 = 2.5)

```
61    --Q7
62    SELECT claim_flag,ROUND(AVG(VehAge),2)
63    FROM Auto_insurance
64    GROUP BY claim_flag;
65
66
67
68
```

|   | claim_flag | ROUND(AVG(VehAge),2) |
|---|------------|----------------------|
| 1 | 0 | 7.07 |
| 2 | 1 | 6.5 |

**INFERENCE- The vehicle age of people who claim their insurance is less when compared with the vehicle age of people who don't claim their insurance,this can be taken in this way also people care a lot about newly bought bikes.**

7.2. Now if we calculate the average Vehicle Age for those who claimed and group them by Area, what do you see in the summary? Any particular pattern you see in the data? (1.5+1=2.5)

Ashwath J

```
65    --PART 2
66    SELECT Area,AVG(VehAge) AS veh_age
67    FROM Auto_insurance
68    WHERE claim_flag = 1
69    GROUP BY Area
70    ORDER BY Area;|
71
```

|   | Area | veh_age |
|---|------|---------|
| 1 | A | 7.43407162078245 |
| 2 | B | 6.97988980716253 |
| 3 | C | 6.44025224454895 |
| 4 | D | 6.49011657374557 |
| 5 | E | 6.09772478070175 |

**INFERENCE- The vehicle age is continuously decreasing.**

Q8). If we calculate the average vehicle age by exposure bucket(as mentioned above), we see an interesting trend between those who claimed vs those who didn't. What is that?

Ashwath J

```
72    --Q8
73    SELECT bucket_list,ROUND(AVG(VehAge),2) AS veh_age_claimed
74    FROM (
75      SELECT * ,
76      CASE
77          WHEN Exposure  between 0 and 0.25 THEN 'E1'
78          WHEN Exposure  between 0.26 and 0.5 THEN 'E2'
79          WHEN Exposure  between 0.51 and 0.75 THEN 'E3'
80          ELSE 'E4'
81      END bucket_list
82      FROM Auto_insurance)
83    WHERE claim_flag = 1
84    GROUP BY bucket_list
85    ORDER BY bucket_list;
86
```

|   | bucket_list | veh_age_claimed |
|---|---|---|
| 1 | E1 | 4.9 |
| 2 | E2 | 6.22 |
| 3 | E3 | 6.18 |
| 4 | E4 | 7.42 |

|   | bucket_list | veh_age_notclaimed |
|---|---|---|
| 1 | E1 | 6.37 |
| 2 | E2 | 6.72 |
| 3 | E3 | 6.27 |
| 4 | E4 | 8.31 |

**INFERENCE- There is no much increase in the veh_age between claimed and not claimed except in E1 bucket.**

Q9)
9.1. Create a Claim_Ct flag on the ClaimNb field as below, and take average of the BonusMalus by Claim_Ct. (2)

```
101        --Q9
102      ⊟SELECT Claim_Ct,ROUND(AVG(BonusMalus),2) AS avg_bonusmalus FROM (
103        │ SELECT * ,
104      ⊟CASE
105        │     WHEN ClaimNb = 1 THEN '1 Claim'
106      ⊟ │     WHEN ClaimNb > 1 THEN 'MT1 Claim'
107        │     ELSE 'No Claim'
108      ├END Claim_Ct
109      └FROM Auto_insurance)
110        GROUP BY Claim_Ct;
111        │
```

| | Claim_Ct | avg_bonusmalus |
|---|----------|----------------|
| 1 | 1 Claim | 62.84 |
| 2 | MT1 Claim | 67.55 |
| 3 | No Claim | 59.59 |

9.2. What is the inference from the summary? (1)

**INFERENCE - The average fine is more for people who have claimed more than once.**

Q10) Using the same Claim_Ct logic created above, if we aggregate the Density column (take average) by Claim_Ct, what inference can we make from the summary data?(4) Note: 2.5 Marks for SQL and 1.5 for inference.

Ashwath J

```
112    --Q10
113    □SELECT Claim_Ct,ROUND(AVG(Density),2) AS avg_density FROM (
114      SELECT * ,
115    □CASE
116        WHEN ClaimNb = 1 THEN '1 Claim'
117    □    WHEN ClaimNb > 1 THEN 'MT1 Claim'
118        ELSE 'No Claim'
119    ├END Claim_Ct
120    └FROM Auto_insurance)
121     GROUP BY Claim_Ct;
122
```

| Claim_Ct | avg_density |
|---|---|
| 1 1 Claim | 1947.32 |
| 2 MT1 Claim | 2297.45 |
| 3 No Claim | 1783.21 |

**INFERENCE - More than one claim are mostly done in high dense cities/areas.A simple analogy is people in metro cities claim more than once than normal tier 2 and tier 3 places.**

Q11) Which Vehicle Brand & Vehicle Gas combination have the highest number of Average Claims (use ClaimNb field for aggregation)? (2)

Ashwath J

```
123     --Q11
124     SELECT VehBrand,VehGas,round(AVG(ClaimNb),3) AS avg_claim
125     FROM Auto_insurance
126     GROUP BY VehGas,VehBrand
127     ORDER BY avg_claim DESC;
128
```

| | VehBrand | VehGas | avg_claim |
|---|---|---|---|
| 1 | B12 | Regular | 0.064 |
| 2 | B5 | Regular | 0.059 |
| 3 | B13 | Diesel | 0.057 |
| 4 | B5 | Diesel | 0.057 |
| 5 | B1 | Regular | 0.054 |

**The B12 Regular** model has the highest average claim among the other models.

Q12)  List the Top 5 Regions & Exposure[use the buckets created above] Combination from Claim Rate's perspective. Use claim_flag to calculate the claim rate. (3)

Ashwath J

```
129    --Q12
130    SELECT Region,bucket_list,SUM(claim_flag) AS claims
131    FROM(
132    SELECT *,CASE
133        WHEN Exposure  between 0 and 0.25 THEN 'E1'
134        WHEN Exposure  between 0.26 and 0.5 THEN 'E2'
135        WHEN Exposure  between 0.51 and 0.75 THEN 'E3'
136        ELSE 'E4'
137    END bucket_list
138    FROM Auto_insurance)
139    GROUP BY bucket_list,Region
140    ORDER BY claims DESC
141    LIMIT 5;
142
```

|   | Region | bucket_list | claims |
|---|--------|-------------|--------|
| 1 | R24 | E4 | 5225 |
| 2 | R82 | E4 | 2258 |
| 3 | R53 | E4 | 1592 |
| 4 | R93 | E4 | 1268 |
| 5 | R24 | E3 | 1221 |

Q13)

Ashwath J

13.1. Are there any cases of illegal driving i.e. underaged folks driving and committing accidents? (1)

```
143     --Q13
144     SELECT IDpol,ClaimNb
145     FROM(
146       SELECT * ,
147     CASE
148       WHEN DrivAge=18 then 'Beginner'
149       WHEN DrivAge BETWEEN 19 AND 30 then 'Junior'
150       WHEN DrivAge BETWEEN 31 AND 45 then 'Middle Age'
151       WHEN DrivAge BETWEEN 46 AND 60 then 'Mid-Senior'
152       WHEN DrivAge >60 then 'Senior'
153       ELSE 'Illegal'
154       END AS Age_flag
155       FROM Auto_insurance)
156       WHERE DrivAge = 'Illegal';
```

Result: 0 rows returned in 512ms

**INFERENCE - No accident cases are committed by driver less than 18 years**

13.2. Create a bucket on DrivAge and then take the average of BonusMalus by this Age Group Category. WHat do you infer from the summary? (2.5+1.5 = 4) Note: DrivAge=18 then 1-Beginner, DrivAge<=30 then 2-Junior, DrivAge<=45 then 3- Middle Age, DrivAge<=60 then 4-Mid-Senior, DrivAge>60 then 5-Senior.

Ashwath J

```
143     --Q13
144     SELECT Age_flag,ROUND(Avg(BonusMalus),2) as avg_bonusmalus
145     FROM(
146     SELECT * ,
147     CASE
148     WHEN DrivAge=18 then 'Beginner'
149     WHEN DrivAge BETWEEN 19 AND 30 then 'Junior'
150     WHEN DrivAge BETWEEN 31 AND 45 then 'Middle Age'
151     WHEN DrivAge BETWEEN 46 AND 60 then 'Mid-Senior'
152     WHEN DrivAge >60 then 'Senior'
153     ELSE 'Illegal'
154     END AS Age_flag
155     FROM Auto_insurance)
156     GROUP BY Age_flag
157     ORDER BY avg_bonusmalus DESC;
```

|   | Age_flag   | avg_bonusmalus |
|---|------------|----------------|
| 1 | Beginner   | 93.01          |
| 2 | Junior     | 79.43          |
| 3 | Middle Age | 59.41          |
| 4 | Mid-Senior | 53.95          |
| 5 | Senior     | 52.8           |

**INFERENCE - The inexperienced drivers i.e, drivers starting their driving career are the people paying most of the fines.**

**CONCEPTUAL QUESTIONS:**

Q14. Mention one major difference between unique constraint and primary key? (2)
(A)Unique constraints can have null values also included but the primary key should be unique and not null.

Q15. If there are 5 records in table A and 10 records in table B and we cross-join these two tables, how many records will be there in the result set? (2)
(A)Cross join is a cartesian product of the number of rows in each table,So therefore the result set will be having 5*10 i.e,50 rows.

Q16. What is the difference between inner join and left outer join? (2)
(A)Inner join does not return null values but left outer join returns null values if there is no exact match on the right table.

Ashwath J

Q17. Consider a scenario where Table A has 5 records and Table B has 5 records. Now while inner joining Table A and Table B, there is one duplicate on the joining column in Table B (i.e. Table A has 5 unique records, but Table B has 4 unique values and one redundant value). What will be the record count of the output? (2)
(A)The output table will have 4 sets which are in common with both the tables.

Q18. What is the difference between WHERE clause and HAVING clause? (2)
(A) WHERE clause is used to filter the table values before GROUP BY but Having clause
Is used to filter the output after the GROUP BY clause.

# THE END