Fall
2023

# MGMT 58200:
# SQL Project for
# Market Data Forecast

GROUP_3
PRATHYUSHA REDDY MIDUDHULA
VENKATA SAI TEJA GANGUMALLA
GOUTHAM KUMAR VEMASANI
HARISH DATTA CHITNENI
MAX MCTIGUE
SANJAY KATYAL

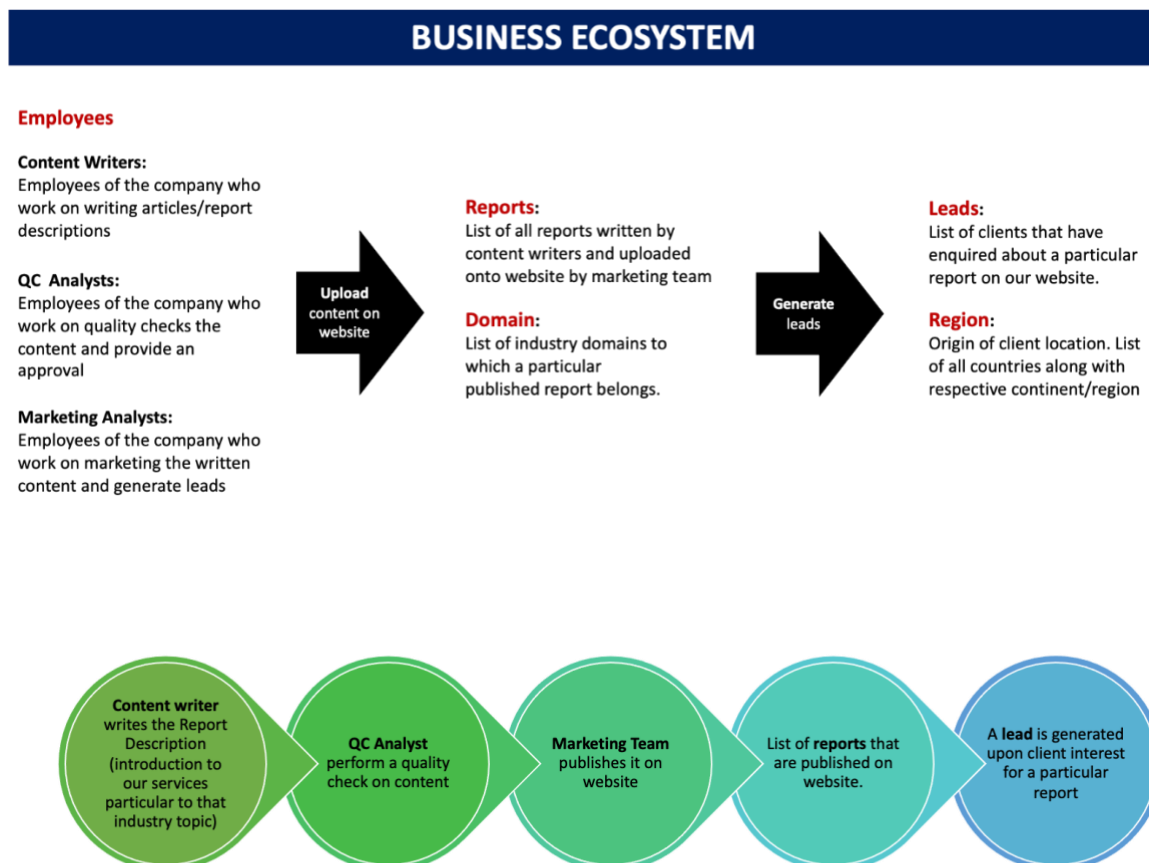PURDUE UNIVERSITY | West Lafayette

## Table of Contents

## About the Company:

- Market Data Forecast (MDF) is a provider of syndicated and custom-made market research, business intelligence and consulting services on gamut of sectors across the globe.

- Primarily, the company offers pre-published market research reports for clients. In addition to this, the company only offers custom market/business research services in-line with client objectives.

- The company has been associated with many Fortune500 listed companies and offers its services spanning across 12+ industry domains and 60+ sub-categories.

## Business Ecosystem and Process:

## Business Case objectives:

→ Migrate the existing Excel-based data system to a robust SQL database to enhance data integrity, security, and accessibility.

→ Analyze the total volume of content created by the company to assess its impact on overall company objectives.

→ Develop an SQL process to calculate the monthly payment for content creators based on a per-word payment model.

→ Implement an SQL-based performance evaluation system to assess the quality and effectiveness of content created by individual content writers. Provide recommendations for improvements.

→ Evaluate the performance of each marketing team member using SQL-based metrics and KPIs to determine their contributions to the company's marketing efforts.

→ Monthly Performance Growth Tracking

→ Establish SQL-based tracking and reporting mechanisms to monitor the month-over-month growth rate in performance metrics for both content creation and marketing activities.

## Attributes Description:

| Entity: Employees |
| --- |
| **EmployeeID:** Unique identifier for each employee. |
| **Name:** The name of the employee. |
| **Title:** The job title or position of the employee. |
| **Pricer_Per_Word:** The pricing rate per word for the employee's services. |
| **Commision_rate:** The commission rate for the employee. |
| **Commision_limit:** The maximum commission limit for the employee. |

| Entity: Reports |
| --- |
| **ReportID:** Unique identifier for each report. |
| **AddedDate:** Date when the report title was added by marketing team member. |
| **WordCount:** Minimum number of words of content required for the report. |
| **ContentWriterID:** Identifier of the content writer who authored the report. |
| **WrittenDate:** Date when the report was written by content writer. |
| **QC_Date:** Date when the report underwent quality control (QC). |
| **QCManagerID:** Identifier of the QC manager responsible for QC. |
| **QC_Status:** Status of the quality control process. |
| **UpdatedDate:** Date when the report was last updated. |

| Entity: Domain |
| --- |
| **DomainID:** Unique identifier for each domain. |
| **Domain:** The name or description of the domain. |

| Entity: Region |
| --- |
| **Country:** Name of a country. |
| **Continent:** The continent to which the country belongs. |

| Entity: Leads |
| --- |
| **id:** Unique identifier for each lead. |
| **report_id:** Identifier linking the lead to a specific report. |
| **created_at:** Date and time when the lead was created. |
| **country:** The country associated with the lead. |

## Normalization Analysis:

Market Data Forecast has provided the data in three tables, first is the Report_Table and second is the Leads_Table and the third is employees table, which are all in two normal form:

- There are no multi-valued or composite attributes
- They **do not have any partial dependencies**, i.e. all data is related to one primary key. "ReportID" in Report_Table and "id" in Leads_Table



**Converting the 2 Normal Form to 3 Normal Form:**

There are transitive dependencies, in both the tables, Report_Table and the Leads_Table.

In **Reports_Table**:

ReportID → DomainID → Domain.

ReportID → ContentWriterID → ContentWriter.

ReportID → QCManagerID → QCManager.

ReportID → ContentWriterID → price_per_word.

ReportID → DomainID → MarketingAnalyst.



In **Leads_Table**:

Id → country → continent

Id → ReportID → Domain

**To remove these transitive dependencies, we moved few attributes to new tables.**

## ER Diagram:



## ER Schema:

## SQL QUERIES AND OUTPUT

**Objective 1: Find the amount of money to be paid to each Content writer for the current month till date.**

```
2      -- Q1 - latest month payment of CW
3  •   select name, cur_month, round(ppw*tot_words) as total_pay from
4  ⊖   (
5      SELECT employees.Name, monthname(updateddate) as cur_month,
6       employees.Pricer_Per_Word as ppw,  sum(reports.WordCount) AS tot_words
7      FROM reports
8      join employees
9      on employees.EmployeeID = reports.ContentWriterID
10     where month(updateddate) = month(curdate())-1
11     group by employees.Name, monthname(updateddate), employees.Pricer_Per_Word)a
12     order by 3 desc;
13
```

**Output:**

| name | cur_month | total_pay |
|------|-----------|-----------|
| Susmriti | September | 2143 |
| Hima Bindu | September | 528 |
| Abhinav | September | 293 |
| Sanya | September | 261 |
| Azim | September | 203 |
| Anuja | September | 190 |
| Akshita | September | 185 |
| Ankita | September | 174 |
| Akanksha | September | 166 |
| Revathi | September | 148 |
| sanket | September | 110 |

**Business Application**: This gives us the amount we need to pay each of our employees for the current month and this can help us in keeping track of payments to be made.

**Objective 2 - Total amount to be spent by the company in the current month.**

```
15    -- Q2 - Net monthly expenditure of company on CW's
16 •  select cur_month, sum(round(ppw*tot_words)) as total_pay from
17  ⊖  (
18    SELECT employees.Name, monthname(WrittenDate) as cur_month,
19    employees.Pricer_Per_Word as ppw,  sum(reports.WordCount) AS tot_words
20    FROM reports
21    join employees
22    on employees.EmployeeID = reports.ContentWriterID
23    group by employees.Name, monthname(WrittenDate), employees.Pricer_Per_Word)a
24    group by cur_month;
```

**Output:**

| cur_month | total_pay |
|---|---|
| November | 277 |
| December | 388 |
| June | 666 |
| July | 2124 |
| August | 2744 |
| April | 540 |
| May | 1387 |
| September | 1433 |
| October | 152 |
| February | 63 |
| March | 33 |

**Business Application**: This gives us the amount we need to pay in total to all our content writers in the current month and this can help in gauging how much budget we need for the month.

**Objective 3 - Commission earned by Content writers in the latest month based on threshold.**

```
1       -- Q3 Commision earned by cw in the latest month based on threshold
2    •  select employeeid, month1,
3    ⊖  case
4       when
5    ⊖  c_l < (case when L_m_commission is null then 0 else L_m_commission end) then 0
6       when
7       cumulate_commission > c_l and c_l > (case when L_m_commission is null then 0 else L_m_commission end)
8    ⊖          then c_l - (case when L_m_commission is null then 0 else L_m_commission end)
9       when
10      c_l > (case when L_m_commission  is null then 0 else L_m_commission  end)
11   ⊖          then cumulate_commission - (case when L_m_commission is null then 0 else L_m_commission end)
12   └          end as current_month_pay,
13   ⊖  case when L_m_commission < c_l then L_m_commission
14      when L_m_commission is null then 0
15      when L_m_commission > c_l then c_l end
16   └  as commision_earned_until_last_month
17   ⊖          from(
18      select *, lag(cumulate_commission,1) over(partition by employeeid order by month1) as L_m_commission
19      from
20   ⊖  (select * , sum(commision_earned) over(partition by employeeid order by month1) as cumulate_commission
21      from
22   ⊖  (select employeeid, month1, commision_rate* lead_count as commision_earned, c_l from
23   ⊖  (select employeeid, commision_rate, c_l , monthname(created_at) as month1, count(id) as lead_count
24      from
25   ⊖  (select employeeid, reportid, commision_rate, commision_limit as c_l
26      from employees as a
27      join reports as b
28      on a.employeeid = b.contentwriterid
29   ┌  where title = 'Content Writer') as c
30      join leads as d
31      on c.reportid = d.report_id
32      where year(created_at) = 2023
33   ┌  group by employeeid, commision_rate, c_l, month1)e)f)g)h
34   └  where month1 = monthname(date_sub(curdate(), interval 1 month))
35      order by 4 desc;
```

**Output**:

| employeeid | month1 | current_month_pay | commision_earned_until_last_month |
|---|---|---|---|
| 22 | September | 0 | 100 |
| 31 | September | 0 | 100 |
| 36 | September | 0 | 100 |
| 41 | September | 0 | 100 |
| 26 | September | 10 | 85 |
| 28 | September | 5 | 75 |
| 18 | September | 20 | 45 |
| 39 | September | 10 | 35 |
| 16 | September | 5 | 30 |
| 17 | September | 10 | 30 |
| 21 | September | 10 | 20 |

**Business Application**: This gives us the amount of commission earned by our content writers, based on the performance of their articles. This helps in understanding our bonus payment budget

**Objective 4 – Calculating the average number of days it takes to Publish an article after it has passed the Quality check**



**Output**:



**Business Application**: Helps in understanding the amount of time it takes to complete the Quality check process, to help make it more efficient if necessary.

**Objective 5- Arranging the Marketing analysts by the lead count their report has produced**

```
33      -- Q5 Best Marketing Analyst
34 •    select d.employeeid, name, sum(leadcount) as leadcount from
35    ⊖ (select c.EmployeeID, count(a.id) as leadcount
36      from leads as a
37      left join reports as b
38      on a.report_id = b.ReportID
39      left join domain as c
40      on c.DomainID = b.DomainID
41      where year(a.created_at) = year(current_date())
42      group by c.EmployeeID
43      having c.EmployeeID is not NULL)d
44      join
45      employees e
46      on d.employeeid = e.employeeid
47      where title = 'Analyst'
48      group by 1,2;
49
```

**Output**:

| employeeid | name | leadcount |
|---|---|---|
| 1 | Bhavesh | 3580 |
| 2 | Karthik | 538 |
| 3 | Prashanth | 636 |
| 4 | Vamshi | 201 |
| 5 | Madhu | 652 |
| 6 | Manoj | 563 |
| 7 | Jathin | 408 |
| 10 | Bhargav | 2889 |

**Business Application**: To figure out the top performing marketing analysts who have produced the most lead count on their article, so as to know the top performers.

**Objective 6 – What is the proportion of reports in each domain that are generating leads within the first 3 months.**

```
51 •     select distinct domain, article_proportion*100 as `article_proportion in %` from
52  ⊖  (select   domainid,
53  ⊖  count(case when id is not null
54          and updateddate between date_sub(created_at, INTERVAL 90 DAY) and
55          created_at then report_id else null end)/count(report_id) as article_proportion
56      from reports as a
57      left join
58      leads b
59      on a.reportid = b.report_id
60      group by 1
61      )c
62      join
63      domain d
64      on c.domainid = d.domainid
65      order by 2 desc
66      ;
67
```

**Output**:

| domain | article_proportion in % |
|---|---|
| Food and Beverage | 21.4600 |
| Health Care | 18.6700 |
| Energy and Resources | 13.3100 |
| Hospitality & Tourism | 10.4400 |
| Automotive | 8.6700 |
| Information Technology | 4.6100 |
| Aerospace And Defense | 4.4800 |
| Automation and Process Control | 3.3100 |
| Electronics and Semiconductor | 2.3300 |
| Consumer Goods And Services | 1.0400 |

**Business Application**: From this data, we can see what proportion of the articles from each domain are able to generate leads within the first three months, this gives us an idea about the best performing domains.

**Objective 7– To find all the content writers with zero leads in last quarter**

```
68      -- Q7 content writer with zero leads in last quarter
69
70  •   select distinct name, count(id) as leads_cnt from
71  ⊖   (select a.name, b.reportid, b.domainid from
72       employees as a
73       inner join
74       reports as b
75       on a.employeeid = b.contentwriterid
76       where title = 'Content Writer') as c
77       left join
78       (select * from leads where created_at between date_sub(CURDATE(), INTERVAL 90 DAY) AND CURDATE())d
79       on c.reportid = d.report_id
80       group by 1
81       having leads_cnt = 0;
82
```

**Output**:

| name | leads_cnt |
| --- | --- |
| Janice | 0 |
| Natania | 0 |
| sanket | 0 |
| Shreya | 0 |
| Swetha | 0 |
| Teja | 0 |
| Vinod | 0 |

**Business Application**: From this list we can see the lowest performing content writers for the last four months. It can help deciding the future pay per word for the writers based on performance.

**Objective 8 - Percentage change in Marketing Analyst Performance Month-over-month (leads generated in M-1 vs leads generated in M-2)**

```
84        -- Q8 Percentage change in Marketing Analyst Performance Month-over-month (leads generated in M-1 vs leads
85  •   SELECT a.employeeid,
86              ((m2_leadcount - m1_leadcount) * 100 / m1_leadcount) AS percent_lead_change
87      FROM (
88          SELECT c.EmployeeID,
89              COUNT(CASE WHEN MONTH(CURDATE()) - MONTH(created_at) = 1
90              THEN a.id ELSE NULL END) AS m1_leadcount,
91              COUNT(CASE WHEN MONTH(CURDATE()) - MONTH(created_at) = 2
92              THEN a.id ELSE NULL END) AS m2_leadcount
93          FROM leads AS a
94          LEFT JOIN reports AS b ON a.report_id = b.ReportID
95          LEFT JOIN domain AS c ON c.DomainID = b.DomainID
96          WHERE YEAR(a.created_at) = YEAR(CURDATE())
97          AND c.employeeid IN (SELECT DISTINCT employeeid FROM employees WHERE title = "Analyst")
98          GROUP BY c.EmployeeID
99      ) AS a
100     WHERE a.employeeid IS NOT NULL

101     order by 2 desc;
102
103     -- which domain is creating more leads
104  •  select domain.Domain, count(leads.id) as leadcount
105     from leads
106     left join reports
107     on leads.report_id = reports.ReportID
108     left join domain
109     on domain.DomainID = reports.DomainID
110     GROUP by domain
111     order by leadcount DESC;
```

**Output**:

| employeeid | percent_lead_change |
|------------|---------------------|
| 6          | 62.8571             |
| 3          | 41.6667             |
| 1          | 24.5552             |
| 5          | 17.0213             |
| 2          | 11.9048             |
| 10         | 10.0877             |
| 4          | -26.6667            |
| 7          | -29.2683            |

**Business Application**: We can use this as a metric to see how performances are changing month over month, based on the leads that each marketing analyst is generating in a given month.

**Objective 9 – To find the most lead generating Domain**

```
103     -- which domain is creating more leads
104  •  select domain.Domain, count(leads.id) as leadcount
105     from leads
106     left join reports
107     on leads.report_id = reports.ReportID
108     left join domain
109     on domain.DomainID = reports.DomainID
110     GROUP by domain
111     order by leadcount DESC;
112
```

**Output**:

| Domain | leadcount |
|---|---|
| Health Care | 4965 |
| Food and Beverage | 3891 |
| Information Technology | 932 |
| Chemicals & Materials | 771 |
| Agriculture | 595 |
| Aerospace And Defense | 513 |
| Automotive | 427 |
| Energy and Resources | 263 |
| Hospitality & Tourism | 249 |
| Automation and Process Control | 181 |
| Consumer Goods And Services | 96 |

**Business Application**: From this table we can see which domain/industry related reports are generating the most leads and thus we can allot more resources to the top performing domains.

**Objective 10 - Find the performance of marketing team members monthly.**

```
130    -- Q10 marketing team members performance till date per month
131 •  select domain.EmployeeID, month(leads.created_at) as leadmonth, count(leads.id) as leadcount
132    from leads
133    left join reports
134    on leads.report_id = reports.ReportID
135    left join domain
136    on domain.DomainID = reports.DomainID
137    where year(leads.created_at) = year(current_date()) and
138    employeeid in (select employeeid from employees where title = 'Analyst')
139    group by month(leads.created_at), domain.EmployeeID
140    having domain.EmployeeID is not NULL;
141
```
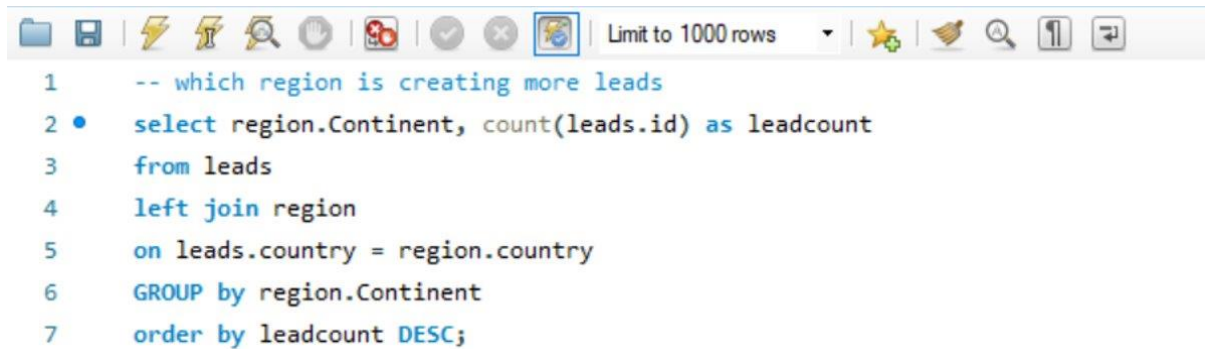
**Output**:

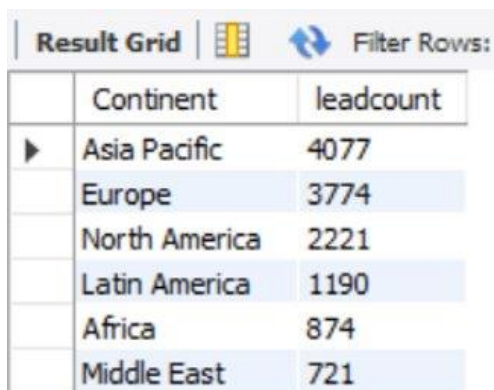| EmployeeID | leadmonth | leadcount |
|---|---|---|
| 1 | 1 | 341 |
| 10 | 1 | 224 |
| 2 | 1 | 55 |
| 3 | 1 | 54 |
| 4 | 1 | 20 |
| 1 | 2 | 399 |
| 10 | 2 | 301 |
| 2 | 2 | 70 |
| 6 | 2 | 50 |
| 3 | 2 | 69 |
| 7 | 2 | 43 |

**Business Application**:  The result of this query helps in understanding the month wise performance of each marketing member, which can help identify the best performers and also the consistency level in performances.

**Objective 11 - To find the continent which generates most number of leads.**



**Output**:



| Continent | leadcount |
|---|---|
| ▶ Asia Pacific | 4077 |
| Europe | 3774 |
| North America | 2221 |
| Latin America | 1190 |
| Africa | 874 |
| Middle East | 721 |

**Business Application**:  The result of this query helps in identifying the region/continent with most number of leads and identify top attractive regions in terms of market potential and plan the marketing activities accordingly.

**Objective 12 – Few summary KPI's that can help the organization.(group of queries)**

```
1     -- Number of leads generated in previous month
2     select count(*) as no_of_leads
3     from leads
4     where month(created_at) = month(now())-1;
5     -- Number of reports written in previous month
6  ●  select count(ReportID) as no_of_reports_written
7     from reports
8     where month(WrittenDate) = month(now())-1;
9     -- Number of reports uploaded in previous month
10 ●  select count(ReportID) as no_of_reports_uploaded
11    from reports
12    where month(UpdatedDate) = month(now())-1;
```

**Output**:

Number of leads generated in last month

| | no_of_leads |
|---|---|
| ▶ | 1433 |

Number of reports written by content writers in last month

| | no_of_reports_written |
|---|---|
| ▶ | 51 |

Number of reports uploaded onto website in the past month

| | no_of_reports_uploaded |
|---|---|
| ▶ | 332 |

## Recommendations:

**Data Integrity and Validation:**

→ Implement data validation checks to minimize manual data errors and ensure data accuracy.

**Structured Storage:**

→ Store all data in a structured format, adhering to the Third Normal Form (3NF), for easy reference in the future.

**Scalability and Performance:**

→ Design the database for computational efficiency, allowing it to handle large data volumes with minimal system lags.

→ Utilize appropriate indexing and partitioning strategies for performance optimization.

**Data Volume:**

→ Recognize that Excel is not suitable for storing extensive data sets due to limitations; an SQL database is better equipped for handling large amounts of data.

→ By following these recommendations, you can maintain data accuracy, structure, and performance, which Excel may struggle to handle for extensive data storage.