

Database Management

Homework 3

B11705044

Yen-Hung, Chiang

1

(a)

根據題目的敘述查詢到的員工編號有四個，下面為查詢的 SQL 程式碼。查詢結果如 Table 1 所表示。

```
1  SELECT
2      E.Ssn,
3      E.Bdate,
4      COUNT(DISTINCT DEP.Dependent_name) AS dependent_count
5  FROM
6      EMPLOYEE E
7  JOIN
8      DEPARTMENT D ON E.Dno = D.Dnumber
9  JOIN
10     DEPT_LOCATIONS DL ON D.Dnumber = DL.Dnumber
11 LEFT JOIN
12     DEPENDENT DEP ON E.Ssn = DEP.Essn
13 GROUP BY
14     E.Ssn, E.Bdate
15 HAVING
16     COUNT(DISTINCT DL.Dlocation) > 1
17 ORDER BY
18     E.Ssn;
```

Ssn	Bdate	dependent_count
123456789	1965-01-09	3
333445555	1955-12-08	3
453453453	1972-07-31	0
666884444	1962-09-15	0

Table 1: Employee Information with Birthdate and Dependent Count

(b)

根據題目的敘述有兩個部門符合條件，下面為查詢所需的 SQL 程式碼。查詢的結果如 Table 2 所表示

```

1  SELECT
2      D.Dnumber AS Department_Number,
3      D.Dname AS Department_Name,
4      FLOOR(AVG(FLOOR(DATE_PART('year', AGE('2024-10-01', E.Bdate))))) AS Average_Age
5  FROM
6      EMPLOYEE E
7  JOIN
8      DEPARTMENT D ON E.Dno = D.Dnumber
9  GROUP BY
10     D.Dnumber, D.Dname
11  HAVING
12     COUNT(E.Ssn) > 1;

```

DNumber	DName	Average_Age
4	Administration	64
5	Research	60

Table 2: Average Age of Employees in Departments with More Than Two Employees

(c)

根據題目的敘述要查詢的條件為參與人數第二多的專案，由於 (e) 小題裡面有提到如果有某個人在某個專案的工時為 NULL 就忽略那筆資料，最後查詢出來的結果有三個專案並列第二名，下面是查詢所需的 SQL 程式碼。查詢結果如 Table 3 所表示。

```

1  WITH ProjectParticipation AS (
2      SELECT
3          P.Pnumber AS Project_Number,
4          P.Pname AS Project_Name,

```

```

5      COUNT(W.Essn) AS Participant_Count,
6      AVG(W.Hours) AS Average_Hours
7  FROM
8      PROJECT P
9  LEFT JOIN
10     WORKS_ON W ON P.Pnumber = W.Pno
11 WHERE
12     W.Hours IS NOT NULL
13 GROUP BY
14     P.Pnumber, P.Pname
15 ),
16 RankedProjects AS (
17     SELECT
18         Project_Number,
19         Project_Name,
20         Average_Hours,
21         Participant_Count,
22         DENSE_RANK() OVER (ORDER BY Participant_Count DESC) AS Participation_Rank
23     FROM
24         ProjectParticipation
25 )
26
27 SELECT
28     Project_Number,
29     Project_Name,
30     Average_Hours
31 FROM
32     RankedProjects
33 WHERE
34     Participation_Rank = 2;

```

Project_Number	Project_Name	Average_Hours
1	ProductX	26.25
3	ProductZ	25.00
20	Reorganization	12.50

Table 3: Average Hours in Projects with Second Highest Participation

(d)

符合條件的員工總共有五位，下面為查詢所需的 SQL 程式碼。查詢結果如 Table 4 所表示。

```

1  WITH ProjectParticipation AS (
2      SELECT
3          P.Pnumber AS Project_Number,
4          P.Pname AS Project_Name,
5          COUNT(W.Essn) AS Participant_Count,
6          AVG(W.Hours) AS Average_Hours
7      FROM
8          PROJECT P
9      LEFT JOIN
10         WORKS_ON W ON P.Pnumber = W.Pno
11     WHERE
12         W.Hours IS NOT NULL
13     GROUP BY
14         P.Pnumber, P.Pname
15 ),
16 RankedProjects AS (
17     SELECT
18         Project_Number,
19         Project_Name,
20         Average_Hours,
21         Participant_Count,
22         DENSE_RANK() OVER (ORDER BY Participant_Count DESC) AS Participation_Rank
23     FROM
24         ProjectParticipation
25 )
26
27 SELECT
28     E.Ssn,
29     E.Lname,
30     CASE
31         WHEN E.Salary >= 40000 THEN 'yes'
32         ELSE 'no'
33     END AS Salary_Above_40k,
34     D.Mgr_ssn,
35     M.Lname AS Manager_Lname
36 FROM
37     EMPLOYEE E
38 JOIN
39     WORKS_ON W ON E.Ssn = W.Essn
40 JOIN
41     PROJECT P ON W.Pno = P.Pnumber

```

```

42 JOIN
43     DEPARTMENT D ON E.Dno = D.Dnumber
44 JOIN
45     EMPLOYEE M ON D.Mgr_ssn = M.Ssn
46 WHERE
47     P.Pnumber IN (
48         SELECT Project_Number
49         FROM RankedProjects
50         WHERE Participation_Rank = 2
51     )
52 AND
53     W.Hours IS NOT NULL
54 GROUP BY
55     E.Ssn, E.Lname, E.Salary, D.Mgr_ssn, M.Lname
56 ORDER BY
57     E.Ssn;

```

Ssn	Lname	Salary_Above_40k	Mgr_ssn	Manager_Lname
123456789	Smith	no	333445555	Wong
333445555	Wong	yes	333445555	Wong
453453453	English	no	333445555	Wong
666884444	Narayan	no	333445555	Wong
987654321	Wallace	yes	987654321	Wallace

Table 4: Employees in Projects with Second Highest Participation

(e)

每位員工的參與的專案數、每週時數總和，以及其參與的專案跨幾個地點如 Table 5 所示，可以看到 Ssn 為”888665555”的這位員工並沒有參與任何專案，所以他的其餘欄位都是 0。下面是查詢所需的 SQL 程式碼。

```

1  SELECT
2      E.Ssn,
3      COALESCE(COUNT(DISTINCT W.Pno), 0) AS Project_Count,
4      COALESCE(SUM(W.Hours), 0) AS Total_Hours,
5      COALESCE(COUNT(DISTINCT P.Plocation), 0) AS Location_Count
6  FROM
7      EMPLOYEE E
8  LEFT JOIN
9      WORKS_ON W ON E.Ssn = W.Essn AND W.Hours IS NOT NULL
10 LEFT JOIN

```

```

11     PROJECT P ON W.Pno = P.Pnumber
12 GROUP BY
13     E.Ssn
14 ORDER BY
15     E.Ssn;

```

Ssn	Project_Count	Total_Hours	Location_Count
123456789	2	40.00	2
333445555	4	40.00	3
453453453	2	40.00	2
666884444	1	40.00	1
888665555	0	0	0
987654321	2	35.00	2
987987987	2	40.00	1
999887777	2	40.00	1

Table 5: Number of Project Participations for every Employee

(f)

所有沒有下屬的員工如 Table 6 所示，總共有五位員工沒有下屬。下面為查詢所需的 SQL 程式碼。

```

1  SELECT
2      E.Ssn,
3      COALESCE(COUNT(DISTINCT W.Pno), 0) AS Project_Count
4  FROM
5      EMPLOYEE E
6  LEFT JOIN
7      EMPLOYEE S ON E.Ssn = S.Super_ssn
8  LEFT JOIN
9      WORKS_ON W ON E.Ssn = W.Essn
10 WHERE
11     S.Ssn IS NULL
12 GROUP BY
13     E.Ssn
14 ORDER BY
15     E.Ssn;

```

Ssn	Project_Count
123456789	2
453453453	2
666884444	1
987987987	2
999887777	2

Table 6: Project Count of Employees Without Subordinates

2

(a)

根據題目的要求得到的查詢結果如 Table 7 所示。下面為這題的 Python 程式碼。

```

1  import pandas as pd
2  import duckdb
3  import psycopg2
4
5  with open('password.txt', 'r') as file:
6      password = file.read().strip()
7
8  connection = psycopg2.connect(
9      dbname='onlinelearning',
10     user='postgres',
11     password=password,
12     host='localhost',
13     port='5432'
14 )
15
16 answers_df = pd.read_sql('SELECT * FROM Answers;', connection)
17 subscriptions_df = pd.read_sql('SELECT * FROM Subscriptions;', connection)
18 state_changes_df = pd.read_sql('SELECT * FROM StateChanges;', connection)
19 user_missions_df = pd.read_sql('SELECT * FROM UserMissions;', connection)
20
21 duck_connection = duckdb.connect()
22
23 duck_connection.register('answers', answers_df)
24 duck_connection.register('subscriptions', subscriptions_df)
25 duck_connection.register('state_changes', state_changes_df)
26 duck_connection.register('user_missions', user_missions_df)

```

```

27
28 query = """
29 SELECT
30     a.AnswerID,
31     a.UserID,
32     a.QuestionID,
33     a.MissionID,
34     a.IsCorrect,
35     a.CostTime,
36     a.CreatedAt,
37     s.EndedAt
38 FROM
39     answers a
40 JOIN
41     subscriptions s ON a.UserID = s.SubscriberID
42 WHERE
43     a.CreatedAt >= '2021-05-01'
44 And
45     a.CreatedAt > s.EndedAt;
46 """
47
48 result_df = duck_connection.execute(query).fetchdf()
49
50 print(result_df.head(5))
51 print(f'Total answers logs: {len(result_df)}')
52
53 connection.close()
54 duck_connection.close()

```

AnswerID	UserID	QuestionID	MissionID	IsCorrect	CostTime	CreatedAT	EndedAt
a3455948	u0004445	q0236673	NaN	1	5	2021-05-01 00:00:07	2021-03-02 23:59:59
a3455952	u0004445	q0236085	NaN	1	4	2021-05-01 00:00:18	2021-03-02 23:59:59
a3455954	u0004445	q0236032	NaN	1	8	2021-05-01 00:00:30	2021-03-02 23:59:59
a3455955	u0004445	q0260118	NaN	1	3	2021-05-01 00:00:39	2021-03-02 23:59:59
a3455958	u0004445	q0234539	NaN	1	3	2021-05-01 00:00:48	2021-03-02 23:59:59

Table 7: First Five Answer Logs After Subscription End Date

(b)

這題的程式碼直接把 (a) 小題得到的結果 `result_df` 直接拿來用，並且把作答紀錄中的平均每題所花時間和總答對率化成一張散佈圖，結果如 Figure 1 所示。可以看到大部分的點都集中在左上方，而平均每題所花時間的中位數為 11.27，總答對率的中位數為 0.83。下面為這題的程式碼。


```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 summary_df = result_df.groupby('userid').agg(
5     AvgTimePerQuestion=('costtime', 'mean'),
6     TotalCorrectRate=('incorrect', 'mean')
7 ).reset_index()
8
9 median_time = summary_df['AvgTimePerQuestion'].median()
10 median_correct_rate = summary_df['TotalCorrectRate'].median()
11
12 plt.figure(figsize=(10, 6), dpi=1000)
13 sns.scatterplot(
14     data=summary_df,
15     x='AvgTimePerQuestion',
16     y='TotalCorrectRate',
17     s=50,
18     color='b'
19 )
20
21 plt.axvline(median_time, color='r', linestyle='--', label=f'Median Time
22     ↳ ({median_time:.2f} sec)')
23 plt.axhline(median_correct_rate, color='g', linestyle='--', label=f'Median Correct
24     ↳ Rate ({median_correct_rate:.2f})')
25
26 plt.xlabel('Avg Time per Question (seconds)')
27 plt.ylabel('Total Correct Rate')
28
29 plt.xlim(0, 50)
30 plt.legend()
31 plt.show()

```

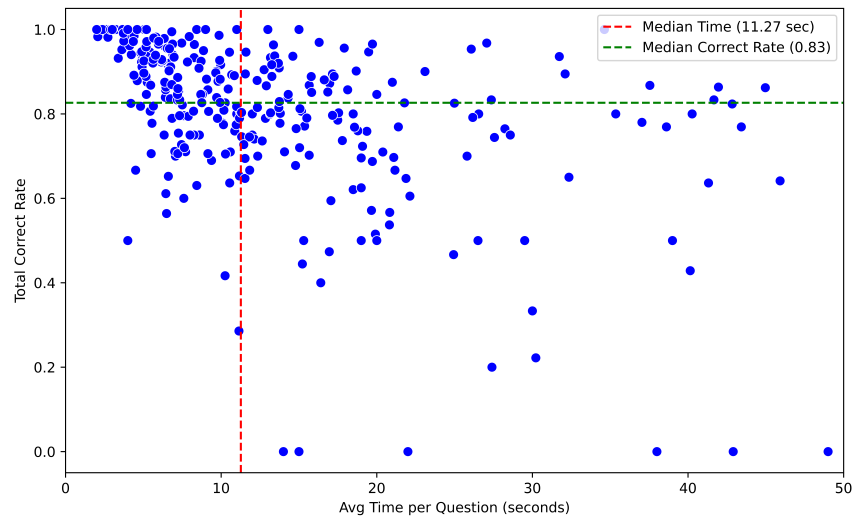


Figure 1: Avg Time per Question v.s. Total Correct Rate

(c)

下面有兩段程式碼分別是任務 i 和任務 ii 的程式碼。

下面這一段程式碼按照任務 i 的敘述依序執行。

```

1  import pandas as pd
2  import time
3
4  answers_df = pd.read_csv('./OnlineLearning_csv/answers.csv')
5  subscriptions_df = pd.read_csv('./OnlineLearning_csv/subscriptions.csv')
6
7  start_time = time.time()
8
9  merged_df = pd.merge(answers_df, subscriptions_df, left_on='UserID',
    ↳ right_on='SubscriberID')
10
11  filtered_df = merged_df[
12      (merged_df['CreatedAt_x'] >= '2021-05-01') &
13      (merged_df['CreatedAt_x'] > merged_df['EndedAt'])
14  ]
15
16  selected_columns = filtered_df[['AnswerID', 'UserID', 'QuestionID', 'MissionID',
    ↳ 'IsCorrect', 'CostTime', 'CreatedAt_x', 'EndedAt']]
17

```

```

18 selected_columns.rename(columns={'CreatedAt_x': 'CreatedAt'}, inplace=True)
19
20 print(selected_columns.head())
21
22 end_time = time.time()
23
24 print(f"Time taken: {end_time - start_time} seconds")

```

而下面這一段程式碼則按照任務 ii 的敘述依序執行。

```

1  import pandas as pd
2  import time
3
4  answers_df = pd.read_csv('./OnlineLearning_csv/answers.csv')
5  subscriptions_df = pd.read_csv('./OnlineLearning_csv/subscriptions.csv')
6
7  start_time = time.time()
8
9  filtered_answers_df = answers_df[answers_df['CreatedAt'] >= '2021-05-01']
10
11 merged_df = pd.merge(filtered_answers_df, subscriptions_df, left_on='UserID',
12   ↪ right_on='SubscriberID')
13
14 final_filtered_df = merged_df[merged_df['CreatedAt_x'] > merged_df['EndedAt']]
15
16 selected_columns = final_filtered_df[['AnswerID', 'UserID', 'QuestionID',
17   ↪ 'MissionID', 'IsCorrect', 'CostTime', 'CreatedAt_x', 'EndedAt']]
18
19 print(selected_columns.head())
20
21 end_time = time.time()
22
23 print(f"Time taken: {end_time - start_time} seconds")

```

執行完上面兩段程式碼之後比較任務 i 和任務 ii 的執行時間，得到的結果如 Table 8 所示，任務 i 的執行時間是任務 ii 的將近五倍。

Task	Time
i	0.8418
ii	0.1497

Table 8

推測任務 i 的執行時間明顯長於任務 ii 的原因為任務 i 會先把整個資料表合併起來再做篩選；而任務 ii 則是先篩選再把表合併起來，透過這種方式可以合併比較少資料。由於合併資料表會花比較多時間，所以表裡面的資料數量很多的話那任務 i 的執行時間就會長很多。

3

(a)

下面是正規化完的 Schema，這個 Schema 沒有 multi-valued attributes，所有非主鍵屬性完全依賴於主鍵，不存在 transitive dependency，每個 determinant 都是 superkey，沒有 multivalued dependency，所以滿足 1 NF, 2 NF, 3 NF, BCNF, 及 4 NF。

Employee_Check

<u>employee_id</u>	<u>work_date</u>	off_type	check_in_time	check_out_time
--------------------	------------------	----------	---------------	----------------

Service

<u>employee_id</u>	<u>work_date</u>	<u>check_out_time</u>	reason
--------------------	------------------	-----------------------	--------

Visit_Customer

<u>employee_id</u>	<u>work_date</u>	<u>check_out_time</u>	<u>customer_id</u>
--------------------	------------------	-----------------------	--------------------

(b)

原本的 relation 就是 2NF 了，因為所有的 nonprime attribute 都依賴於 prime attribute。但不是 3 NF 因為 Fee 依賴 ServiceType，而 ServiceType 又依賴 (EngineerID, CustomerID, ConsultingDate)，所以把它正規化成 3 NF。

Consulting

<u>EngineerID</u>	<u>CustomerID</u>	<u>ConsultingDate</u>	ServiceType
-------------------	-------------------	-----------------------	-------------

Service

<u>ServiceType</u>	Fee
--------------------	-----

這樣就不會有 $X \rightarrow Z$ and $Z \rightarrow Y$ 的關係了。

(c)

設計出的 Schema 如下：

Product				
<u>ProductID</u>	ProductName	Length	Width	Height

Supplier	
<u>SupplierID</u>	SupplierName

Pricing			
<u>ProductID</u>	<u>SupplierID</u>	<u>EffectiveDate</u>	Price

沒有 multi-valued attributes 所以滿足 1 NF。所有非主鍵屬性完全依賴於主鍵，滿足 2 NF。不存在 transitive dependency，滿足 3 NF。每個 determinant 都是 superkey，滿足 BCNF。沒有 mult-valued dependency，滿足 4 NF。

4

(a)

直接使用 SQL 下 query 可以有更好的效率，而且在操作上也會比較靈活，在對性能注重的地方會需要直接用 SQL。但 SQL 的程式碼比較難閱讀，而且重複性高，更有可能會受到 SQL injection 的攻擊。

使用 ORM 的好處是可以防止像是 SQL injection 的攻擊，還可以減少開發者的工作量，另外可以對不同的資料庫操作。缺點是 ORM 的查詢性能沒有 SQL 好，而且學習的成本比較高，對於複雜的查詢靈活度會受限。

(b)

除了 SQL injection 之外還有一種攻擊叫資料庫暴力破解，這種攻擊方式會生成大量的用戶名和密碼組合以猜測用戶的登入資料，反覆嘗試直到成功登入。可以通過限制帳戶登入、多重驗證、監控異常流量等方式進行防範。

(c)

兩者的意思是不太一樣的。「前端與後端切開」的前端意思是負責呈現資料、接收用戶輸入、處理用戶交互的部分；後端則是接收來自前端的請求、處理業務邏輯、從資料庫中檢索或存儲資料的部分。「client 端和 server 端切開」的 client 端主要是在用戶設備上執行的應用程序；server 端主要是在伺服器上執行的應用程序。所以前者主要關注 UI 與業務邏輯的分離，後者主要關注設備與服務的交互。