# LAB – 1

Name: Gandevia Keval Dharmeshbhai
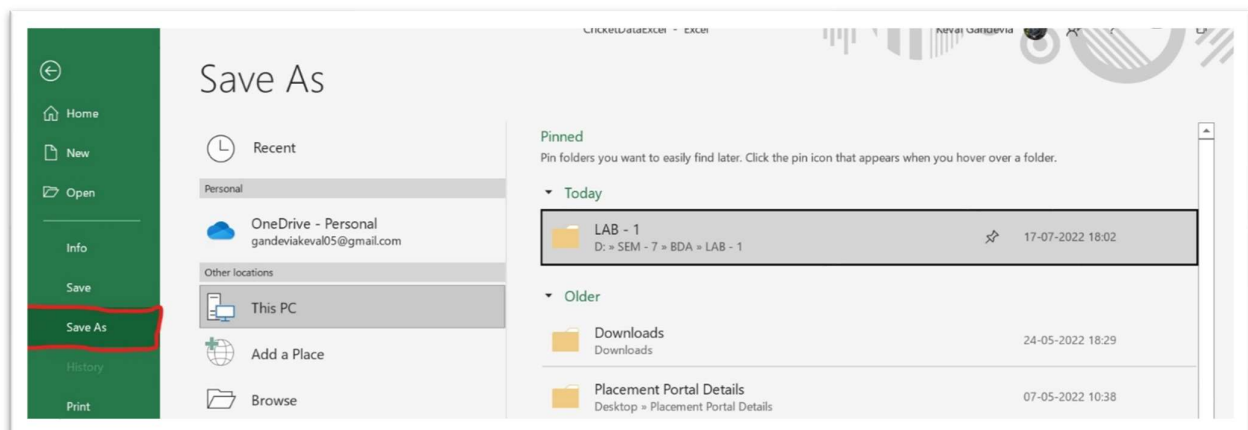
Sem: VII

Roll No: CE046

Subject: Big Data and Analytics
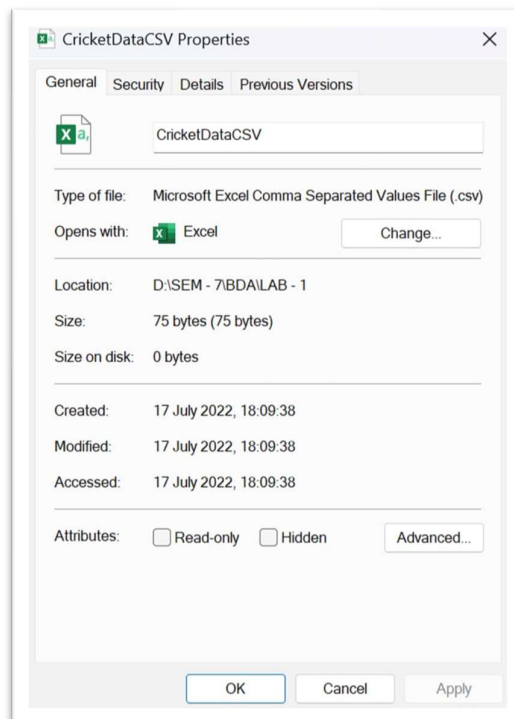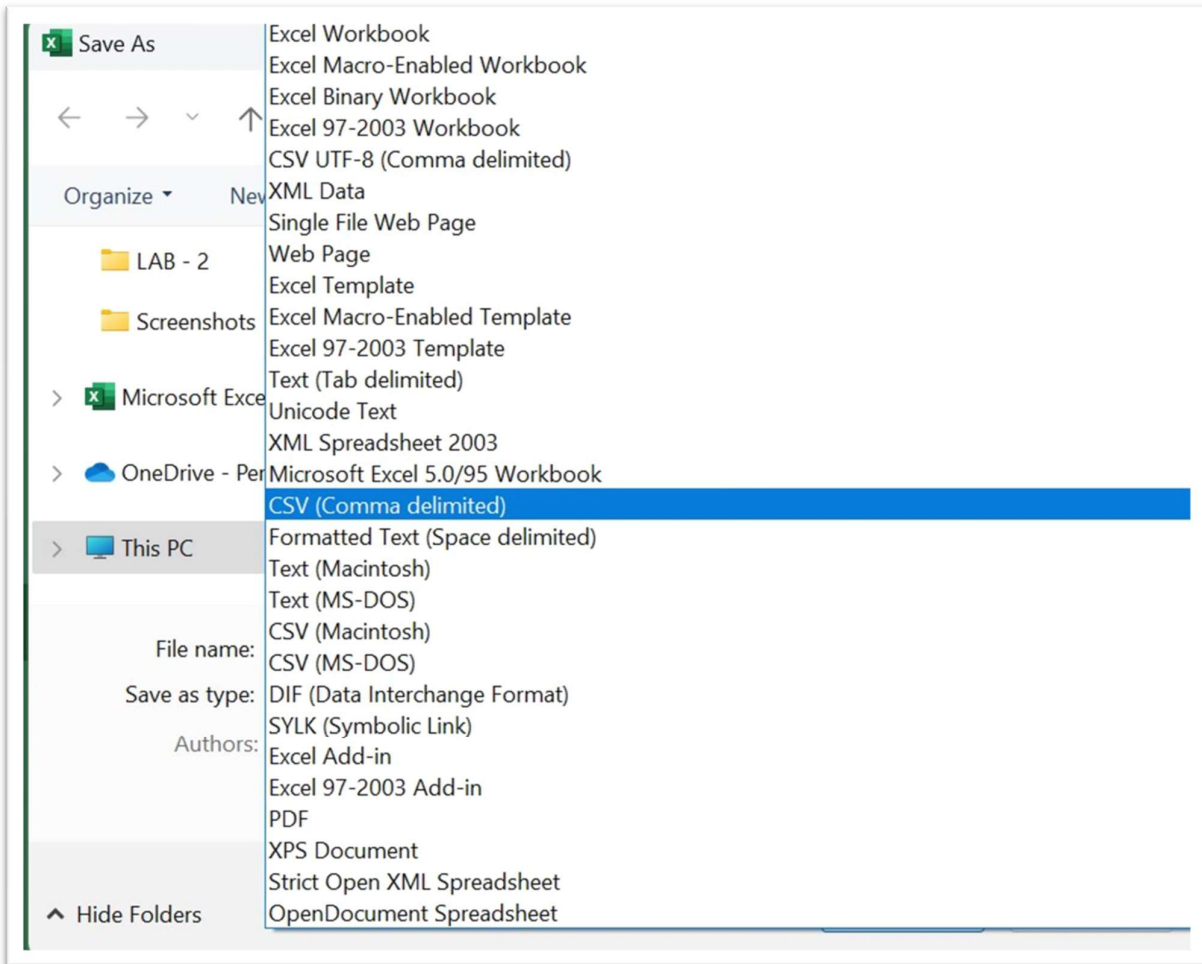
Aim: Recording types of data and various file formats. Identifying data sources. Handling traditionally to start with a small scale.

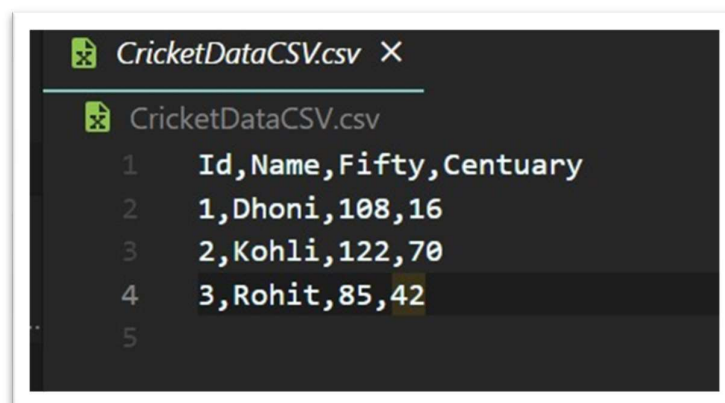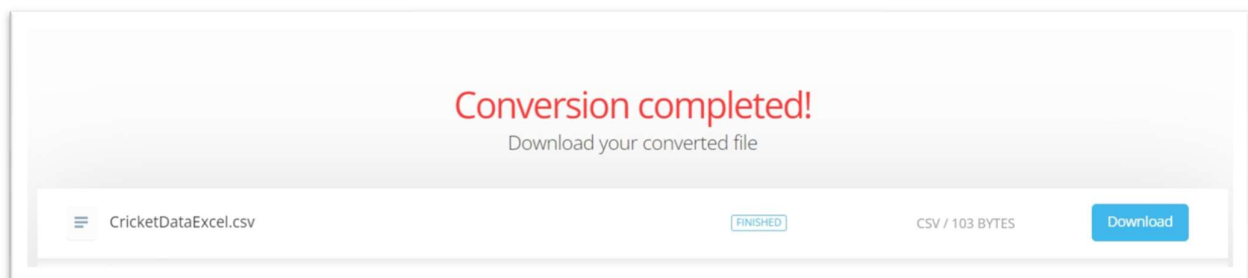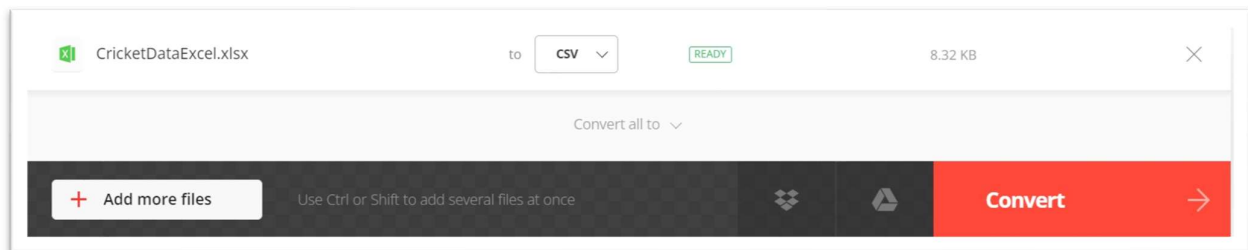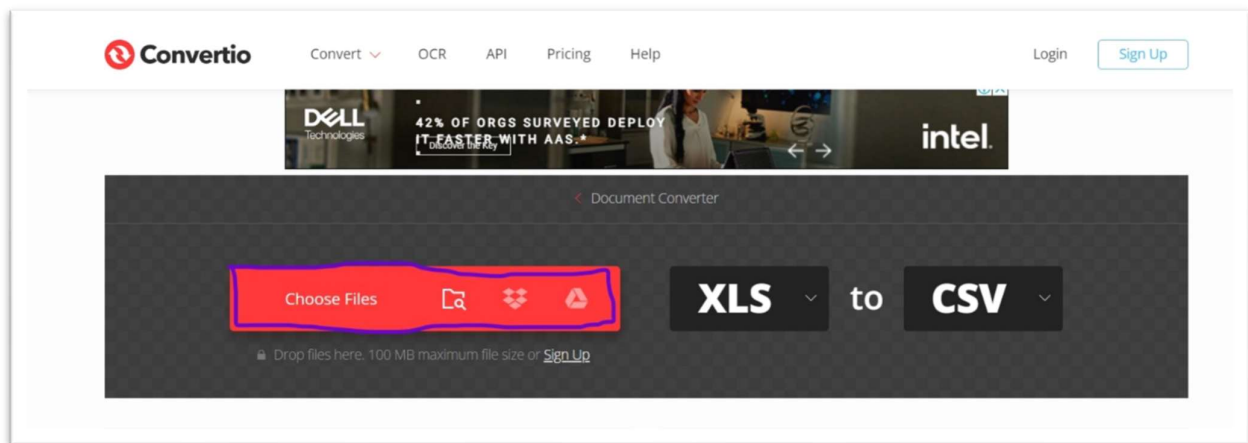## Q. 1: Given the spreadsheet file convert it into a csv.

❖ There are many ways to convert the spreadsheet file into csv.

    I. **Save As spreadsheet file with .csv extension.**

Save As

Organize ▾    New

  📁 LAB - 2

  📁 Screenshots

> 📗 Microsoft Exce

> ☁ OneDrive - Per

> 🖥 This PC

File name:
Save as type:
Authors:

∧ Hide Folders

Excel Workbook
Excel Macro-Enabled Workbook
Excel Binary Workbook
Excel 97-2003 Workbook
CSV UTF-8 (Comma delimited)
XML Data
Single File Web Page
Web Page
Excel Template
Excel Macro-Enabled Template
Excel 97-2003 Template
Text (Tab delimited)
Unicode Text
XML Spreadsheet 2003
Microsoft Excel 5.0/95 Workbook
CSV (Comma delimited)
Formatted Text (Space delimited)
Text (Macintosh)
Text (MS-DOS)
CSV (Macintosh)
CSV (MS-DOS)
DIF (Data Interchange Format)
SYLK (Symbolic Link)
Excel Add-in
Excel 97-2003 Add-in
PDF
XPS Document
Strict Open XML Spreadsheet
OpenDocument Spreadsheet



CricketDataCSV Properties                                    ✕

General  Security  Details  Previous Versions

        CricketDataCSV

Type of file:   Microsoft Excel Comma Separated Values File (.csv)

Opens with:    📗 Excel              Change...

Location:      D:\SEM - 7\BDA\LAB - 1
Size:          75 bytes (75 bytes)
Size on disk:  0 bytes

Created:       17 July 2022, 18:09:38
Modified:      17 July 2022, 18:09:38
Accessed:      17 July 2022, 18:09:38

Attributes:    ☐ Read-only   ☐ Hidden      Advanced...

                    OK        Cancel        Apply

## II. Use online converter.

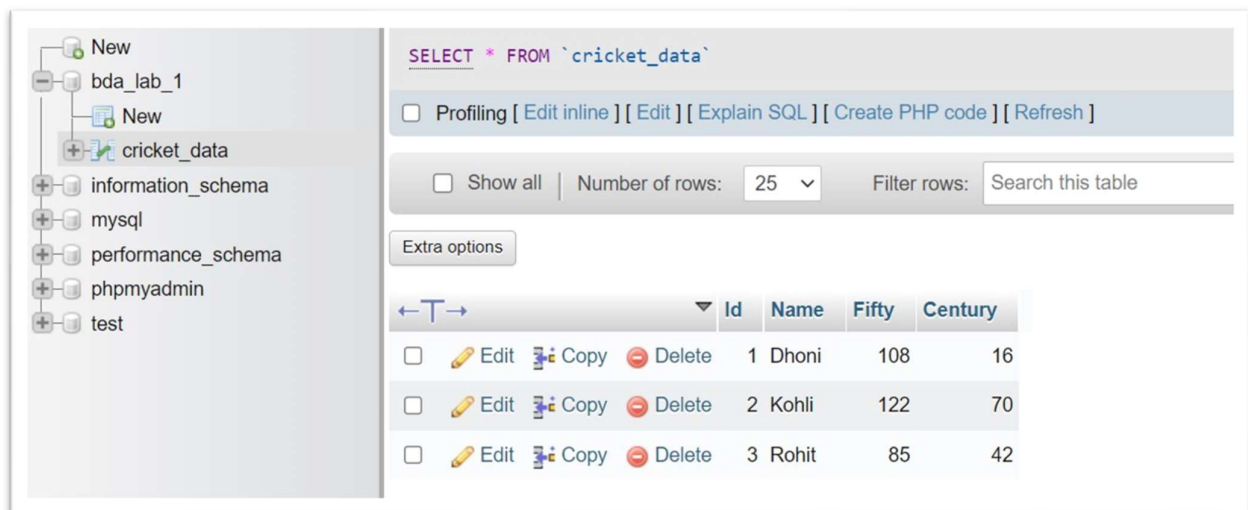# Q. 2: Import csv into MySQL database table.

## ❖ Code:

```
CricketDataCSV.csv        task2.py    ×

task2.py > ...
2    from mysql.connector import Error
3    import mysql.connector as msql
4    import pandas as pd
5    data = pd.read_csv("CricketDataCSV.csv", index_col=False, delimiter=",")
6    print(data.head())
7
8    try:
9        conn = msql.connect(host='localhost', user='root',
10                           database='BDA_LAB_1', password='')
11       if conn.is_connected():
12           cursor = conn.cursor()
13           cursor.execute('''
14               CREATE TABLE cricket_data (
15                   Id INT PRIMARY KEY,
16                   Name VARCHAR(50),
17                   Fifty INT,
18                   Century INT
19               )
20           ''')
21
22           for i, row in data.iterrows():
23               query = "INSERT INTO cricket_data VALUES (%s, %s, %s, %s)"
24               cursor.execute(query, tuple(row))
25               print("Record inserted!!")
26               conn.commit()
27
28   except Error as e:
29       print("Error while connecting to MySQL", e)
30
```

❖ **Output:**

```
[Running] python -u "d:\SEM - 7\BDA\LAB - 1\task2.py"
    Id   Name  Fifty  Centuary
0   1   Dhoni   108        16
1   2   Kohli   122        70
2   3   Rohit    85        42
Record inserted!!
Record inserted!!
Record inserted!!

[Done] exited with code=0 in 0.937 seconds
```

New
bda_lab_1
    New
    cricket_data
information_schema
mysql
performance_schema
phpmyadmin
test

SELECT * FROM `cricket_data`

☐ Profiling [ Edit inline ] [ Edit ] [ Explain SQL ] [ Create PHP code ] [ Refresh ]

☐ Show all | Number of rows: 25 ⌄ | Filter rows: Search this table

Extra options

| ←T→ | | | | Id | Name | Fifty | Century |
|---|---|---|---|---|---|---|---|
| ☐ | 🖉 Edit | ⯗ Copy | ⊖ Delete | 1 | Dhoni | 108 | 16 |
| ☐ | 🖉 Edit | ⯗ Copy | ⊖ Delete | 2 | Kohli | 122 | 70 |
| ☐ | 🖉 Edit | ⯗ Copy | ⊖ Delete | 3 | Rohit | 85 | 42 |

# Q. 3: Write a computer program to read records from database and generate data file.

I.   XML

❖ **Code:**

```
task2.py          task3.py    ×    CricketDataXML.xml

task3.py > ...
1    from msilib.schema import File
2    import mysql.connector as msql
3    from lxml.builder import E
4
5    conn = msql.connect(host='localhost', user='root',
6                        database='BDA_LAB_1', password='')
7    cursor = conn.cursor()
8    cursor.execute("SELECT * FROM cricket_data")
9
10   outfile = open("CricketDataXML.xml", "w")
11   rows = cursor.fetchall()
12   outfile.write('<?xml version="1.0" ?>\n')
13   outfile.write("<CRICKETDATA>\n")
14   for row in rows:
15       outfile.write("  <ROW>\n")
16       outfile.write('    <Id>%s</Id>\n' % row[0])
17       outfile.write('    <Name>%s</Name>\n' % row[1])
18       outfile.write('    <Fifty>%s</Fifty>\n' % row[2])
19       outfile.write('    <Century>%s</Century>\n' % row[3])
20       outfile.write('  </ROW>\n')
21   outfile.write('</CRICKETDATA>\n')
22   outfile.close()
23
```

❖ **Output:**

```xml
CricketDataXML.xml
1    <?xml version="1.0" ?>
2    <CRICKETDATA>
3      <ROW>
4        <Id>1</Id>
5        <Name>Dhoni</Name>
6        <Fifty>108</Fifty>
7        <Century>16</Century>
8      </ROW>
9      <ROW>
10       <Id>2</Id>
11       <Name>Kohli</Name>
12       <Fifty>122</Fifty>
13       <Century>70</Century>
14     </ROW>
15     <ROW>
16       <Id>3</Id>
17       <Name>Rohit</Name>
18       <Fifty>85</Fifty>
19       <Century>42</Century>
20     </ROW>
21   </CRICKETDATA>
```

## II.    JSON

❖ **Code:**

```python
import json
import collections
import mysql.connector as msql

conn = msql.connect(host='localhost', user='root',
                    database='BDA_LAB_1', password='')
cursor = conn.cursor()
cursor.execute("SELECT * FROM cricket_data")
rows = cursor.fetchall()

objects_list = []

for row in rows:
    d = collections.OrderedDict()
    d["Id"] = row[0]
    d["Name"] = row[1]
    d["Fifty"] = row[2]
    d["Century"] = row[3]
    objects_list.append(d)

j = json.dumps(objects_list)

with open("CricketDataJSON.json", "w") as f:
    f.write(j)
```

❖ **Output:**

```
[{"Id": 1, "Name": "Dhoni", "Fifty": 108, "Century": 16}, {"Id": 2, "Name": "Kohli", "Fifty": 122, "Century": 70}, {"Id": 3, "Name": "Rohit", "Fifty": 85, "Century": 42}]
```

# Q. 4: Import XML/JSON file into another database/table.

❖ Creating table in SQL Server.



```
Design        ↑↓      T-SQL
1  ☐CREATE TABLE [dbo].[Cricket_Data]
2   (
3       [Id] INT NOT NULL PRIMARY KEY IDENTITY,
4       [Name] VARCHAR(50) NOT NULL,
5       [Fifty] INT NOT NULL,
6       [Century] INT NULL
7   )
```

❖ Query to convert csv to database table.



```
dbo.Cricket_Data [Data]    dbo.Cricket_Data [Design]    SQLQuery1.sql *  ╕ ×
▷ ▾ ■  ✓  国  ᵌ ᵌ ᵌ  EmployeeDB          ▾  ᵍ  ᵍ ▾ ᵍ  ᵍ
1  ☐SELECT Cricket_Data.*
2   FROM OPENROWSET (BULK 'D:\SEM - 7\BDA\LAB - 1\CricketDataJSON.json', SINGLE_CLOB) as j
3   CROSS APPLY OPENJSON(BulkColumn)
4   WITH( Id int, Name varchar(50), Fifty int, Century int) AS Cricket_Data
```

114 %  ▾   ✓ No issues found                                                          ▶   Ln: 4    Ch:

T-SQL  ↑↓   ☷ Results  ☷ Message

| | Id | Name | Fifty | Century |
|---|---|---|---|---|
| 1 | 1 | Dhoni | 108 | 16 |
| 2 | 2 | Kohli | 122 | 70 |
| 3 | 3 | Rohit | 85 | 42 |

# Q. 5: Export database dump for data migration/archival.

❖ We can easily export any SQL for of database using phpMyAdmin window of Xampp.

## ❖ Output file:

```
bda_lab_1.sql  ✕

bda_lab_1.sql

38    -- Dumping data for table `cricket_data`
39    --
40
41    INSERT INTO `cricket_data` (`Id`, `Name`, `Fifty`, `Century`) VALUES
42    (1, 'Dhoni', 108, 16),
43    (2, 'Kohli', 122, 70),
44    (3, 'Rohit', 85, 42);
45
46    --
47    -- Indexes for dumped tables
48    --
49
50    --
51    -- Indexes for table `cricket_data`
52    --
53    ALTER TABLE `cricket_data`
54      ADD PRIMARY KEY (`Id`);
55    COMMIT;
56
57    /*!40101 SET CHARACTER_SET_CLIENT=@OLD_CHARACTER_SET_CLIENT */;
58    /*!40101 SET CHARACTER_SET_RESULTS=@OLD_CHARACTER_SET_RESULTS */;
59    /*!40101 SET COLLATION_CONNECTION=@OLD_COLLATION_CONNECTION */;
60
```

## Q. 6: Validate/Map data types across different database systems when migrating from one to another.

```
dbo.Cricket_Data [Design]  ⊣ ✕  SQLQuery1.sql *

↥ Update   Script File:  dbo.Table.sql                    ▾

      Name          Data Type    Allow Nulls  Default
⚷ Id              int              ☐
  Name            varchar(50)      ☐
  Fifty           int              ☐
  Century         int              ☐
                                   ☐
```
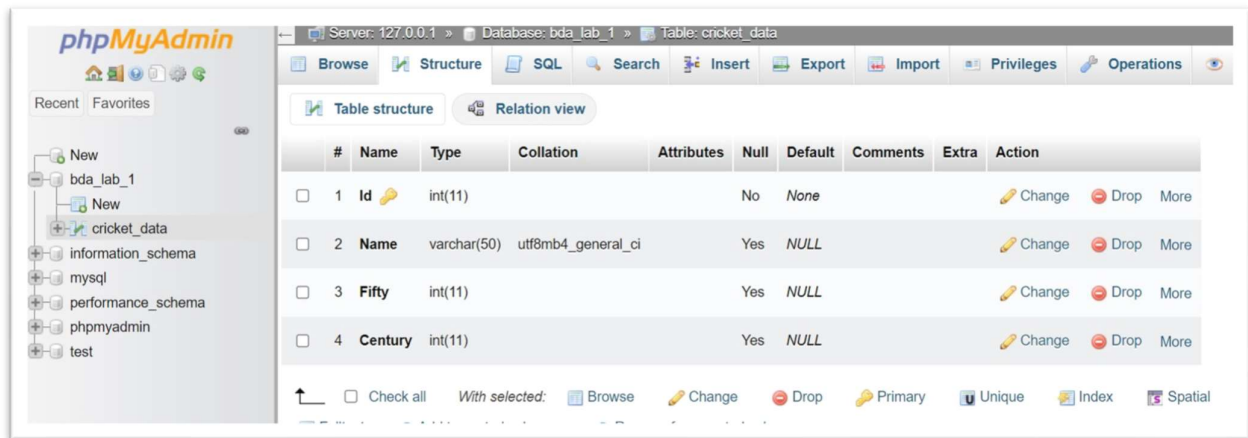
⊿ **Keys** (1)
    <unnamed>  (Primary Key, Clustered: Id)
**Check Constraints** (0)
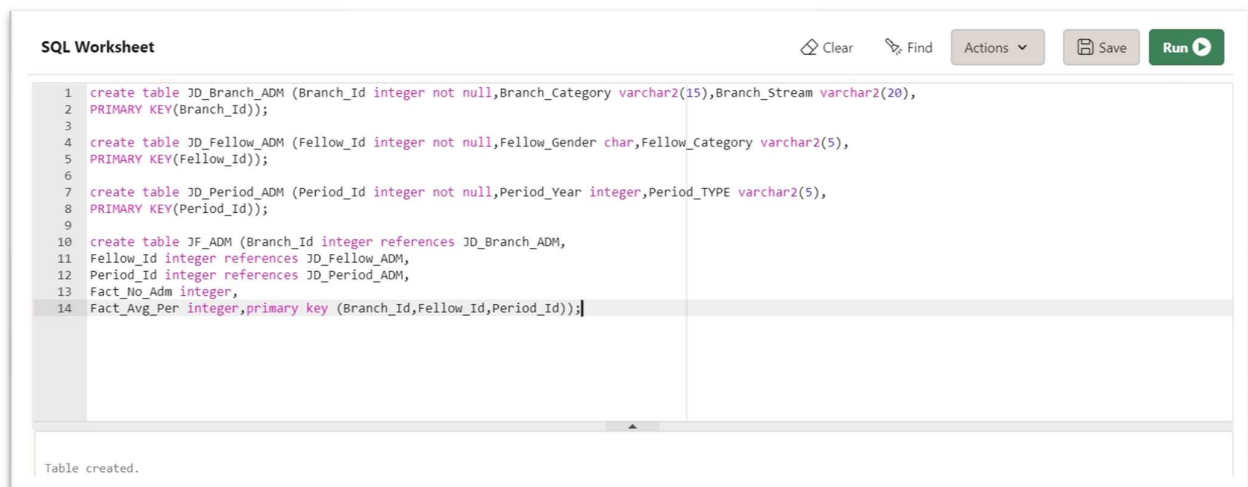**Indexes** (0)
**Foreign Keys** (0)
**Triggers** (0)

## Q. 7: Represent data cube and perform operations. OLTP – Data Warehouse.

❖ **Creating tables:**

❖ **Inserting data into tables:**



```
SQL Worksheet                                  ⟋ Clear   ⟋ Find    Actions ∨    ⊟ Save   Run ▶

16
17  insert into jd_branch_adm values(1,'ENG','CS');
18  insert into jd_branch_adm values(2,'ENG','EC');
19  insert into jd_branch_adm values(3,'MED','MBBS');
20  insert into jd_branch_adm values(4,'MED','DENTAL');
21  insert into jd_fellow_adm values(1,'M','OPEN');
22  insert into jd_fellow_adm values(2,'M','SC');
23  insert into jd_fellow_adm values(3,'F','OPEN');
24  insert into jd_fellow_adm values(4,'F','SC');
25  insert into jd_Period_adm values(1,2004,'1RS');
26  insert into jd_Period_adm values(2,2004,'2RS');
27  insert into jd_Period_adm values(3,2005,'1RS');
28  insert into jd_Period_adm values(4,2005,'2RS');
29  insert into jf_adm values (1,1,1,1000,60);
30  insert into jf_adm values (1,1,2,1100,70);
31  insert into jf_adm values (1,1,3,1200,80);
32  insert into jf_adm values (1,1,4,1300,90);
33  insert into jf_adm values (1,2,1,1400,60);
34  insert into jf_adm values (1,2,2,1500,70);
35  insert into jf_adm values (1,2,3,1600,80);
36  insert into jf_adm values (1,2,4,1700,90);
```

❖ **Representing data as data cube:**



```
SQL Worksheet

93  SELECT BRANCH_ID, FELLOW_ID, PERIOD_ID
94  FROM jf_adm
95  GROUP BY CUBE(BRANCH_ID, FELLOW_ID, PERIOD_ID);
96
97
98
```

| BRANCH_ID | FELLOW_ID | PERIOD_ID |
|---|---|---|
| - | - | - |
| - | - | 1 |
| - | - | 2 |
| - | - | 3 |
| - | - | 4 |
| - | 1 | - |
| - | 1 | 1 |
| - | 1 | 2 |
| - | 1 | 3 |
| - | 1 | 4 |
| - | 2 | - |
| - | 2 | 1 |
| - | 2 | 2 |

**SQL Worksheet**

| | | |
|---|---|---|
| - | 2 | 3 |
| - | 2 | 4 |
| - | 3 | - |
| - | 3 | 1 |
| - | 3 | 2 |
| - | 3 | 3 |
| - | 3 | 4 |
| - | 4 | - |
| - | 4 | 1 |
| - | 4 | 2 |
| - | 4 | 3 |
| - | 4 | 4 |
| 1 | - | - |
| 1 | - | 1 |
| 1 | - | 2 |
| 1 | - | 3 |
| 1 | - | 4 |
| 1 | 1 | - |
| 1 | 1 | 1 |
| 1 | 1 | 2 |
| 1 | 1 | 3 |
| 1 | 1 | 4 |
| 1 | 2 | - |

| | | |
|---|---|---|
| 1 | 2 | 1 |
| 1 | 2 | 2 |
| 1 | 2 | 3 |
| 1 | 2 | 4 |
| 1 | 3 | - |
| 1 | 3 | 1 |
| 1 | 3 | 2 |
| 1 | 3 | 3 |
| 1 | 3 | 4 |
| 1 | 4 | - |
| 1 | 4 | 1 |
| 1 | 4 | 2 |
| 1 | 4 | 3 |
| 1 | 4 | 4 |

Download CSV
Rows 1 - 50. More rows exist.

## ❖ Slice operation:

**SQL Worksheet**

```
98   Select Fact_No_Adm, Fact_Avg_Per
99   from JD_Branch_ADM B,JD_Fellow_ADM F,JD_Period_ADM D,JF_ADM FACT
100  where (
101  B.Branch_Id = FACT.Branch_Id and
102  F.Fellow_Id = FACT.Fellow_Id and
103  D.Period_Id = FACT.Period_Id and
104  F.Fellow_Gender = 'M' and
105  F.Fellow_Category = 'OPEN');
106
107
```

**SQL Worksheet**

| FACT_NO_ADM | FACT_AVG_PER |
|-------------|--------------|
| 1000 | 60 |
| 1100 | 70 |
| 1200 | 80 |
| 1300 | 90 |
| 2600 | 60 |
| 2700 | 70 |
| 2800 | 80 |
| 2900 | 90 |
| 4200 | 60 |
| 4300 | 70 |
| 4400 | 80 |
| 4500 | 90 |
| 5800 | 60 |
| 5900 | 70 |
| 6000 | 80 |
| 6100 | 90 |

Download CSV
16 rows selected.

## ❖ Dice operation:

**SQL Worksheet**

```
107  Select Fact_No_Adm, Fact_Avg_Per
108  from JD_Branch_ADM B,JD_Fellow_ADM F,JD_Period_ADM D,JF_ADM FACT
109  where (
110  B.Branch_Id = FACT.Branch_Id and
111  F.Fellow_Id = FACT.Fellow_Id and
112  D.Period_Id = FACT.Period_Id and
113  (F.Fellow_Id = 1 or F.Fellow_Id = 2) and
114  (B.Branch_Id = 1 or B.Branch_Id = 2) and
115  (D.Period_Id = 1 or D.Period_Id = 2));
116
```

## SQL Worksheet

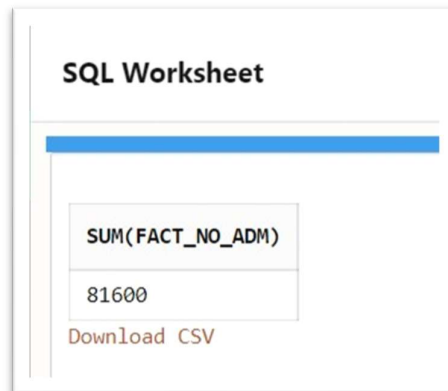| FACT_NO_ADM | FACT_AVG_PER |
|---|---|
| 1000 | 60 |
| 1100 | 70 |
| 1400 | 60 |
| 1500 | 70 |
| 2600 | 60 |
| 2700 | 70 |
| 3000 | 60 |
| 3100 | 70 |

Download CSV
8 rows selected.

## ❖ Roll up operation:

## SQL Worksheet

```
118   Select sum(Fact_No_Adm)
119   from JD_Branch_ADM B,JD_Fellow_ADM F,JD_Period_ADM D,JF_ADM FACT
120   where
121   B.Branch_Id = FACT.Branch_Id and
122   F.Fellow_Id = FACT.Fellow_Id and
123   D.Period_Id = FACT.Period_Id and
124   B.Branch_Category = 'ENG'
125   group by B.Branch_Category;
126
127
```

SQL Worksheet

SUM(FACT_NO_ADM)

81600

Download CSV

## Q. 8: Generate pdf report. / Use any virtualization tool.