



AAiT

ADDIS ABABA INSTITUTE OF TECHNOLOGY

አዲስ አበባ ቴክኖሎጂ ኢንስቲትዩት

ADDIS ABABA UNIVERSITY

አዲስ አበባ ዩኒቨርሲቲ

Department of Information Technology & Scientific Computing (Software Engineering)

Course name: - Big Data Analysis
Lab Report - I

Prepared By: -

Abdulhamid Abdo	ATR/8526/07
Abrham Kindu	ATR/6689/07
Samuel Mussie	ATR/9258/07

Submitted to: - Alazar Alemayehu

Dec 2018

ABSTRACT

Connecting to a database from the R environment can be extremely useful. It allows us direct access to the data without have to first export it from a database and then import it from a csv file or entering it directly into R. This assignment documentation will showcase the processes and steps to set up R to connect to a database. This documentation will show step-by-step instructions and provides the R code on how to query the database once connectivity is established.

Introduction

We all know that R is a very powerful when working with big data. We can manipulate and deal with small datasets that easily fit into your computer's memory. But what if we encounter datasets that are too large for your computer to handle as a whole? In this case, storing the data outside of R and organizing it in a database will become handy. Connecting to the database allows you to retrieve only the chunks needed for the current analysis. R can connect to almost any existing database type. Most common database types have R packages that allow you to connect to them. Examples could be RSQLite, RMySQL, RODBC, etc.). We will be using RODBC to connect to our database and perform CRUD operations on our database.

Problem Description

Performing CRUD (Create, Read, Update, Delete) operations using R to a database.

(A) Design a simple database containing 3 tables with relationships.

- We have created a database called 'BIGDATA' with three tables namely 'USER', 'GRADES', and 'ADDRESS'.
- USER table basic information of a user which has 'id', 'first_name', 'last_name' and 'sex' attributes.
- GRADES table contains information about a user's grade and has 'id', 'user_id', 'course_name' and 'course_grade' attributes. It has 'many to many' relationship with table 'USER' as many grades can be possessed by a user.
- ADDRESS table contains information about a user's address and has 'id', 'user_id', 'house_no', 'region' and 'city' attributes. It has 'one to one' relationship with table 'USER' assuming a user has only one address to live on.

The ER diagram for the database is shown below.

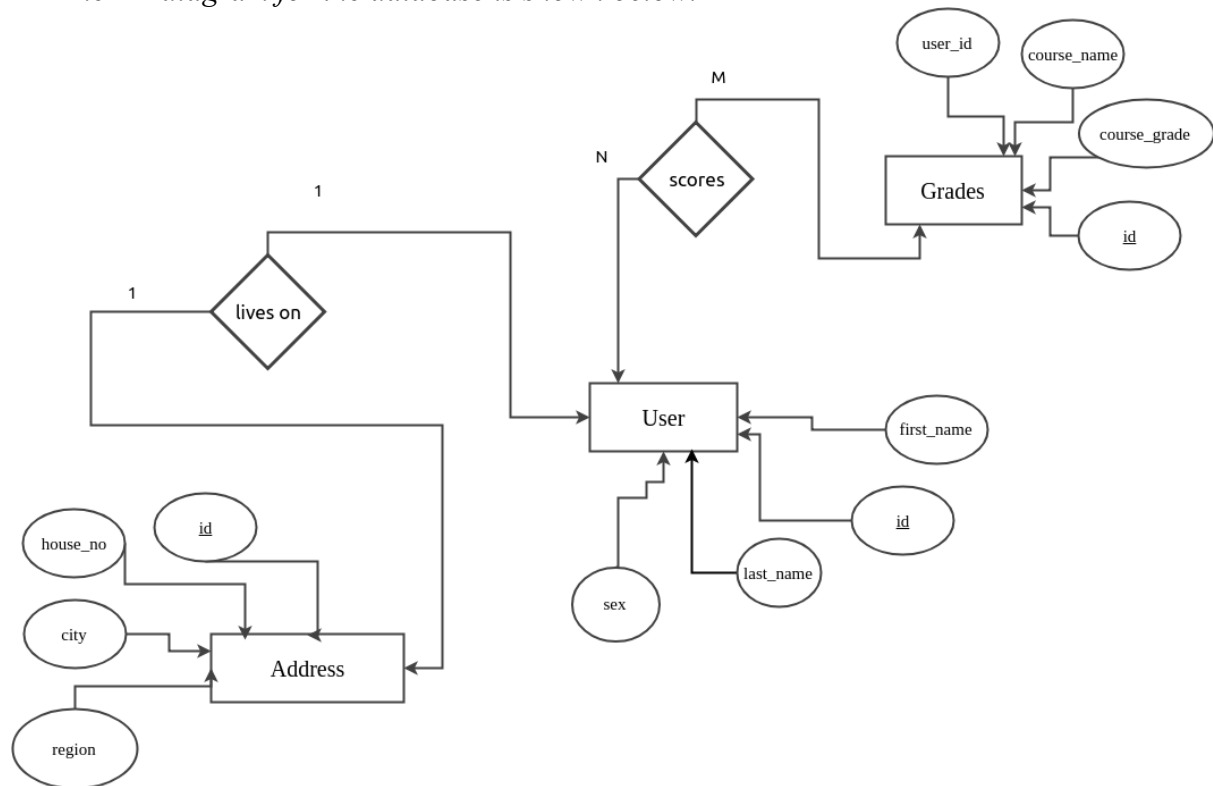


Figure 1 ER diagram for database BIGDATA

(B, C, D & E) Create a table from R program and demonstrate CRUD operations.

- We will be showing step by step processes of creating the tables by using screenshots.

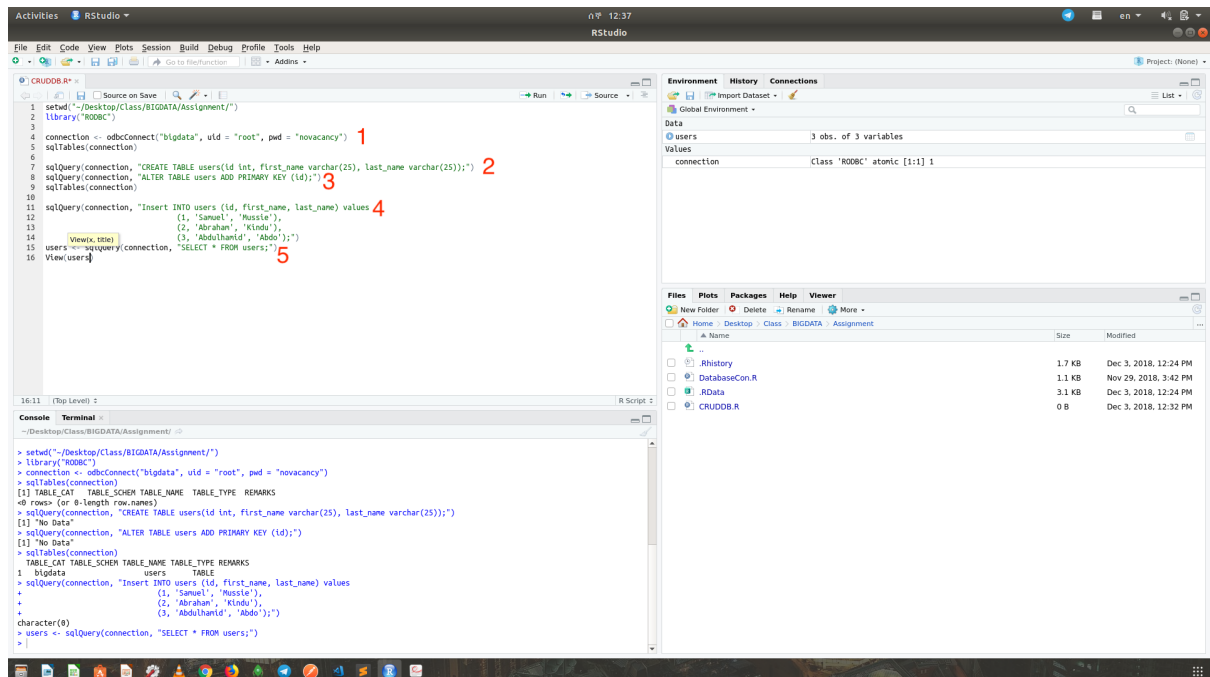


Figure 2 Creating and adding fields to a table using R.

- The above screenshot shows the steps to create a table and add fields to a table using R inside RStudio.
- What the code numbered as 1 on the screenshot does is connecting to a database named 'bigdata' using the provided credentials. The code designated as 2 creates a table named 'users' with the provided attributes. Number 3 adds a 'PRIMARY KEY' 'id' to table 'users'. Number 4 inserts data in to table 'users'. And finally, number 5 assigns the result from the select statement to a variable user in order to view it inside RStudio. The result from executing code number 5 is shown in the figure below.

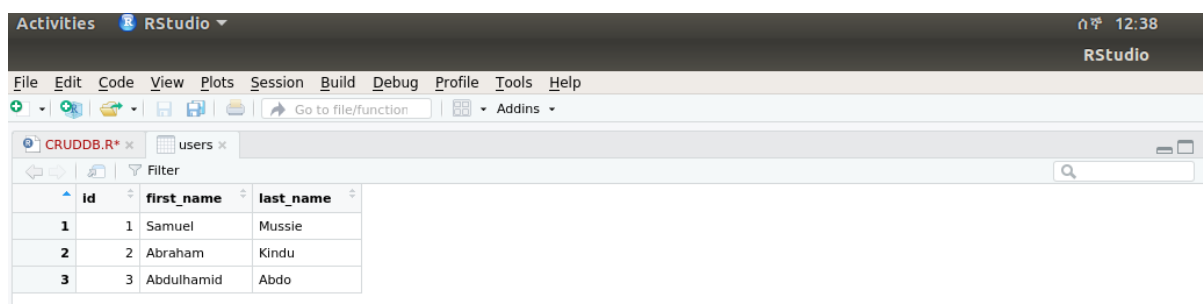


Figure 3 Executing View(users).

- *Executing update operation.*

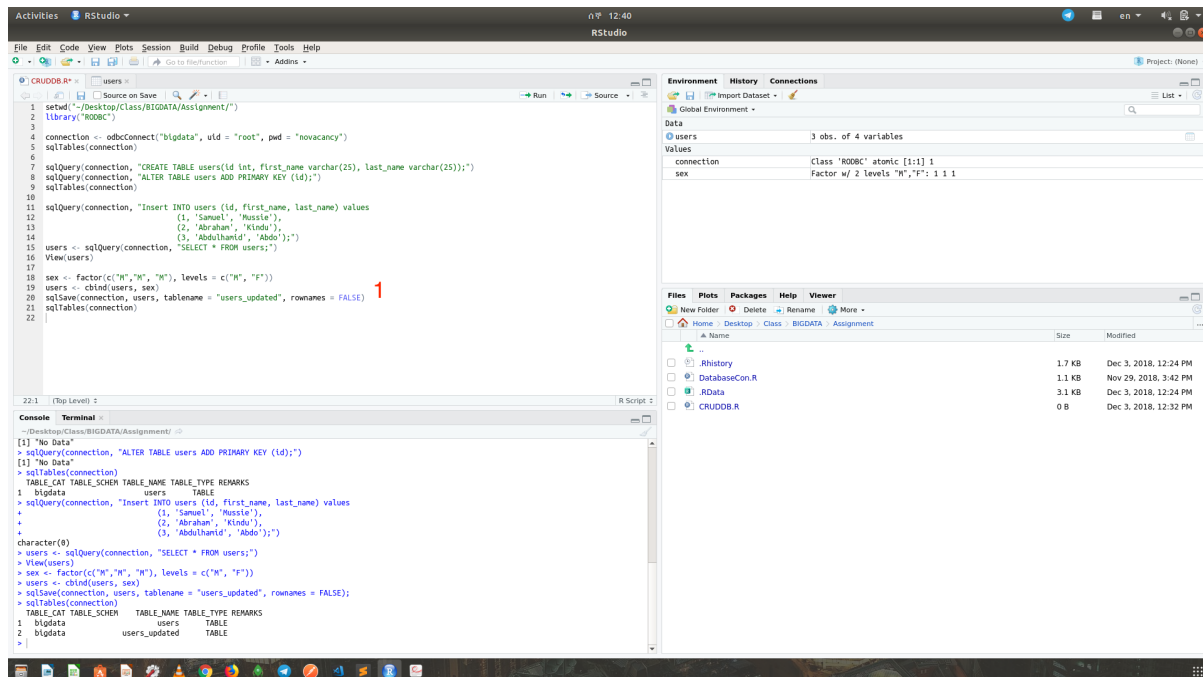


Figure 4 update operation.

- *The codes marked as 1 will factor and add a 'sex' attribute. Then cbind is called to bind 'sex' attribute and then we add it as a column in 'user_updated' table. The contents of the table are shown in the figure below.*

	id	first_name	last_name	sex
1	1	Samuel	Mussie	M
2	2	Abraham	Kindu	M
3	3	Abdulhamid	Abdo	M

Figure 5 Viewing the updated table.

- In similar manners, we create table 'addresses' and 'grades'.
 - Creating table addresses.

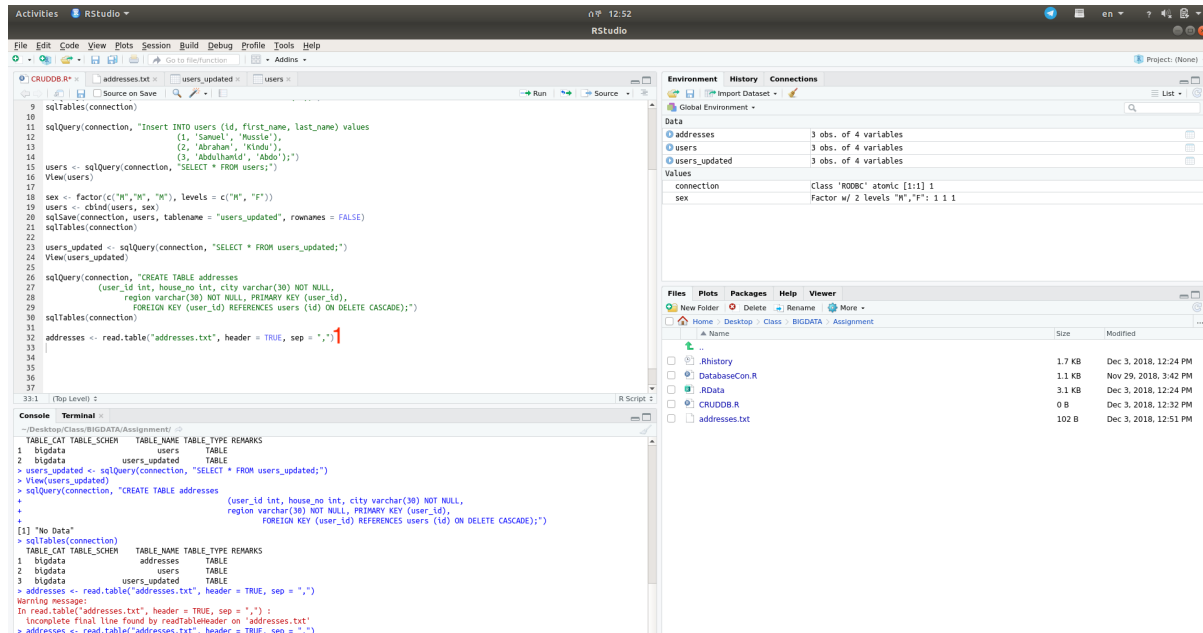


Figure 6 Creating and inserting data to table addresses.

- Code number 1 reads addresses.txt and adds them to table address. The contents of addresses.txt is shown in the figure below.

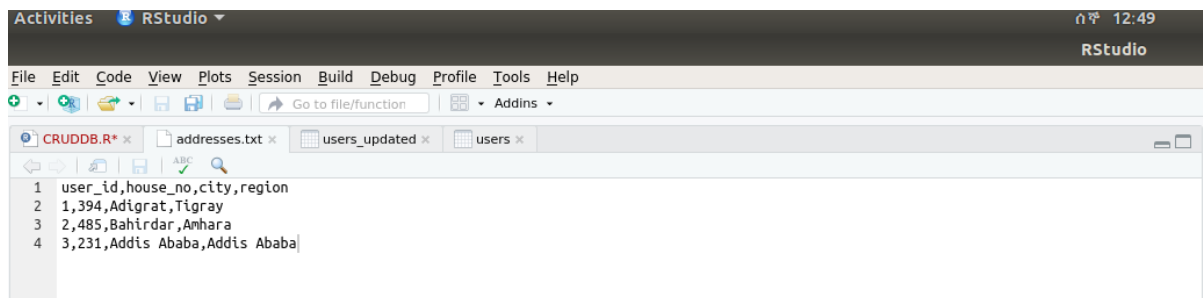


Figure 7 Addresses.txt

- Creating table 'grades'.

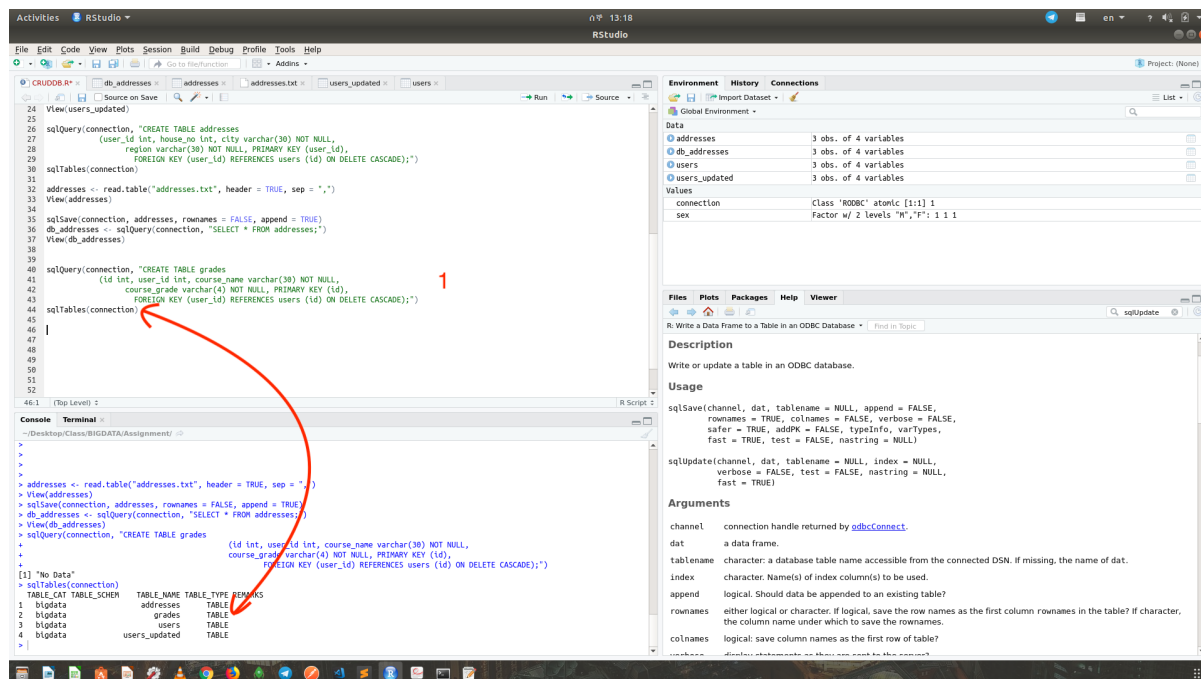


Figure 8 Creating table grades.

- Code numbered as 1 creates a table named grades with the provided attributes. We can see that, after the execution of `sqlTables(connection)`, a new table named grades is added as shown in the console.
- Adding and updating table grades.

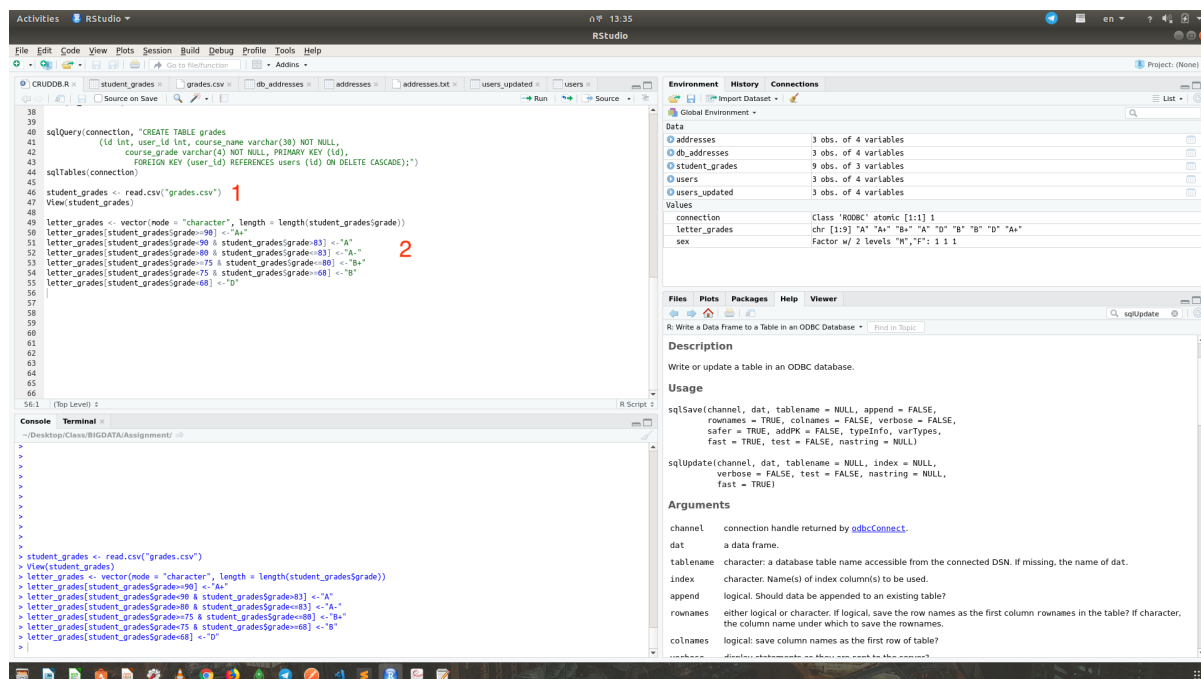
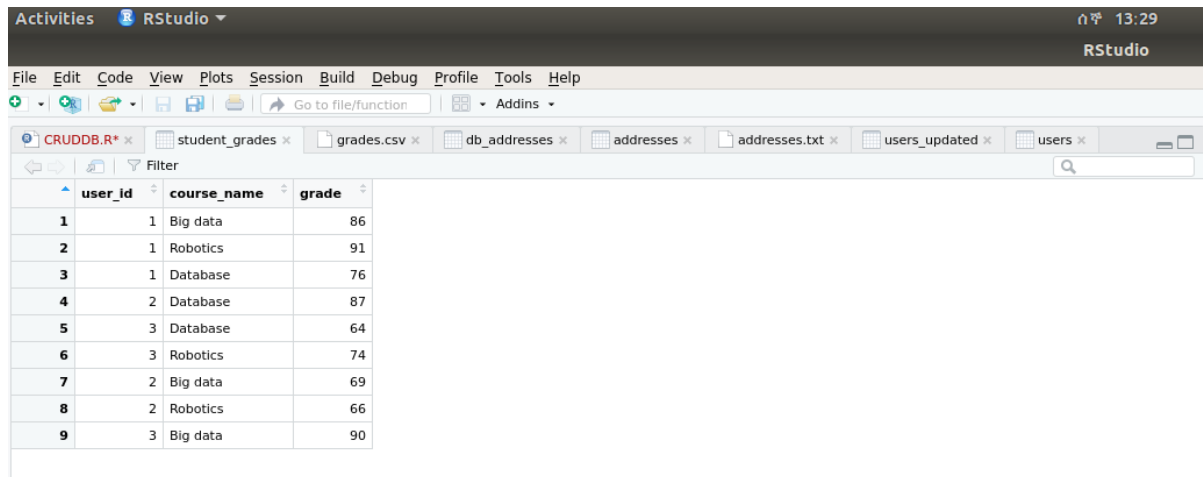


Figure 9 adding and updating table grades.

- Code numbered as 1 reads grades from file grades.csv and assigns it to a variable named `student_grades`. Codes designated as 2 filters student grades and assigns letter grades based on grade values.

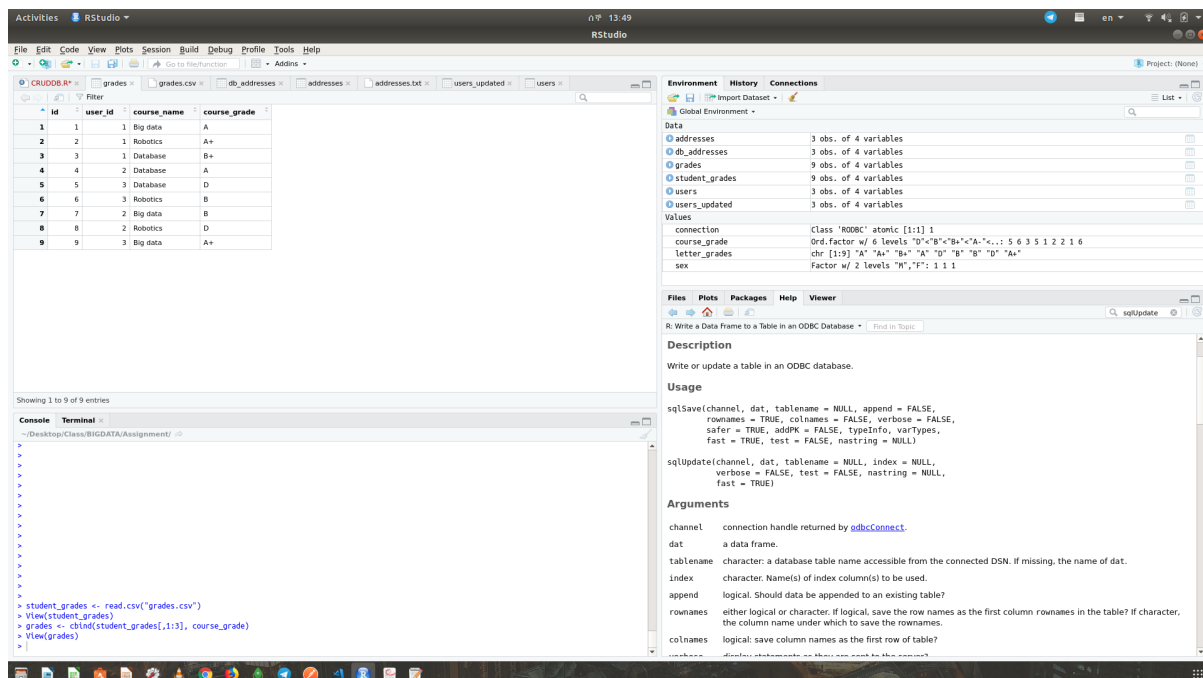
- Contents of student_grades before adding letter grades are shown in the figure below.



	user_id	course_name	grade
1	1	Big data	86
2	1	Robotics	91
3	1	Database	76
4	2	Database	87
5	3	Database	64
6	3	Robotics	74
7	2	Big data	69
8	2	Robotics	66
9	3	Big data	90

Figure 10 Student_grades

- The content of grades table is shown in the figure below.



	id	user_id	course_name	course_grade
1	1	1	Big data	A
2	2	1	Robotics	A+
3	3	1	Database	B+
4	4	2	Database	A
5	5	3	Database	D
6	6	3	Robotics	B
7	7	2	Big data	B
8	8	2	Robotics	D
9	9	3	Big data	A+

```

> student_grades <- read.csv("grades.csv")
> View(student_grades)
> grades <- cbind(student_grades[,1:3], course_grade)
> View(grades)

```

Figure 11 Grades with letter grades

• Update and delete operations

The screenshot shows the RStudio interface with the following components:

- Source Editor:** Contains R code that connects to a database, updates the 'grades' table, and deletes a user.


```

55 letter_grades[student_grades$grade==75 & student_grades$grade==68] <- "B"
56 letter_grades[student_grades$grade==68] <- "D"
57
58 course_grade <- factor(letter_grades, levels = c("D", "B", "B+", "A-", "A", "A+"), ordered = TRUE)
59
60 grades <- chind(student_grades[,1:3], course_grade)
61 View(grades)
62
63 sqlSave(connection, grades, rownames = FALSE, append = TRUE)
64 db_grades <- sqlQuery(connection, "SELECT * FROM grades;")
65 View(db_grades)
66
67 sqlQuery(connection, "UPDATE grades SET course_grade='C' WHERE course_grade ='D';")
68 db_grades <- sqlQuery(connection, "SELECT * FROM grades;")
69 summary(db_grades)
70
71 sqlQuery(connection, "DELETE FROM users WHERE id=1;")
72 users <- sqlQuery(connection, "SELECT * FROM users;")
73 View(users)
74
75 db_addresses <- sqlQuery(connection, "SELECT * FROM addresses;")
76 summary(db_addresses)
77 View(db_addresses)
78
79
80
81
82

```
- Environment Pane:** Shows the state of the database tables:
 - addresses: 3 obs. of 4 variables
 - db_addresses: 2 obs. of 4 variables
 - db_grades: 9 obs. of 4 variables
 - grades: 9 obs. of 4 variables
 - student_grades: 9 obs. of 4 variables
 - users: 2 obs. of 4 variables
 - users_updated: 3 obs. of 4 variables
- Console:** Shows the output of the SQL queries, including the updated 'grades' table and the deleted user.


```

> users <- sqlQuery(connection, "SELECT * FROM users;")
> View(users)
> db_addresses <- sqlQuery(connection, "SELECT * FROM addresses;")
> View(db_addresses)
> summary(db_grades)
  id  user_id  course_name  course_grade
Min.   1 Min.   1 Big data    A :2
1st Qu. 1st Qu. 1 Database    A+ :2
Median   Median 2 Robotics    B :2
Mean     Mean   2 B+ :1
3rd Qu. 3rd Qu. 3 C :2
Max.     Max.   3
> summary(db_addresses)
  user_id  house_no  city  region
Min.   12.00 Min.   2331.0 Addis Ababa:1
1st Qu. 12.25 1st Qu. 2294.5 Bahir Dar :1
Median 12.50 Median 3350.0
Mean    12.50 Mean 3350.0
3rd Qu. 12.75 3rd Qu. 4421.5
Max.    13.00 Max. 4485.0

```

Figure 12 Update and delete operations.

- The code numbered as 1 updates tables grades by looking for course_grades with values 'D' and setting them to value 'C'.
- The code numbered as 2 deletes a user with id 0f 1.
- Updated table and deleted table are shown respectively.

The screenshot shows two side-by-side views of the 'grades' table. The left view shows the original data, and the right view shows the data after the update. Red arrows point from the 'D' values in the 'course_grade' column to the 'C' values in the updated view.

id	user_id	course_name	course_grade
1	1	Big data	A
2	2	Robotics	A+
3	3	Database	B+
4	4	Database	A
5	5	Database	D
6	6	Robotics	B
7	7	Big data	B
8	8	Robotics	D
9	9	Big data	A+

id	user_id	course_name	course_grade
1	1	Big data	A
2	2	Robotics	A+
3	3	Database	B+
4	4	Database	A
5	3	Database	C
6	6	Robotics	B
7	7	Big data	B
8	2	Robotics	C
9	3	Big data	A+

Figure 13 Updated grades table values

The screenshot shows the 'users' table after deleting a user with id 1. The table now only contains two users.

id	first_name	last_name
1	Abraham	Kindu
2	Abdulhamid	Abdo

Figure 14 Deleted User with id 1

- Showing cascading delete from the grades table too. No grade is shown with user id of 1 as user with an id of 1 is deleted his grades will also be automatically deleted.

The screenshot shows the RStudio interface with a database connection. The table view for 'db_grades' is as follows:

id	user_id	course_name	course_grade
1	4	Database	A
2	5	Database	C
3	6	Robotics	B
4	7	Big data	B
5	8	Robotics	C
6	9	Big data	A+

The console shows the following R code and output:

```

3rd Qu.:7 3rd Qu.:3 C :2
Max. :9 Max. :3
> summary(db_addresses)
  user_id house_no city region
Min. :2.00 Min. :231.0 Addis Ababa:1 Addis Ababa:1
1st Qu.:2.25 1st Qu.:294.5 Bahir Dar :1 Amara :1
Median :2.50 Median :358.0
Mean :2.50 Mean :358.0
3rd Qu.:2.75 3rd Qu.:421.5
Max. :3.00 Max. :485.0
> db_grades <- sqlQuery(connection, "SELECT * FROM grades;")
> summary(db_grades)
  id user_id course_name course_grade
Min. :4.00 Min. :2.00 Big data:2 A :1
1st Qu.:15.25 1st Qu.:2.0 Database:2 A+:1
Median :16.50 Median :2.5 Robotics:2 B :2
Mean :6.50 Mean :2.5 C :2
3rd Qu.:17.75 3rd Qu.:3.0
Max. :19.00 Max. :3.0
> View(db_grades)
  
```

The Environment pane shows the following data frames:

- addresses: 3 obs. of 4 variables
- db_addresses: 2 obs. of 4 variables
- db_grades: 6 obs. of 4 variables
- grades: 9 obs. of 4 variables
- student_grades: 9 obs. of 4 variables
- users: 2 obs. of 3 variables
- users_updated: 3 obs. of 4 variables

Figure 15 Demonstrating cascading delete