

데이터 베이스와 R

박찬엽

2017년 6월 27일

목차

- 데이터베이스
 - 데이터베이스란
 - 서버와 클라이언트
 - R과 DB를 연결해주는 DBI
- 데이터 소개
 - 데이터 공유
 - 데이터 원본
 - 데이터 훑어보기
- 클라우드 서비스
 - 클라우드 서비스 소개
 - 구글 클라우드
 - RMySQL 연결

과제 확인

데이터베이스

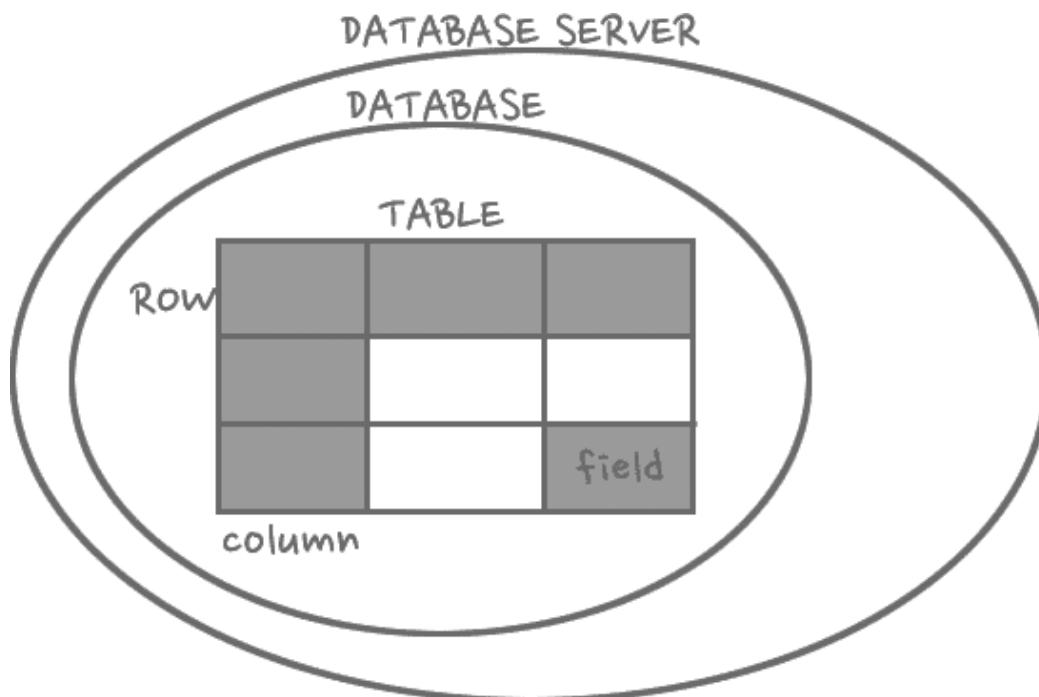
데이터란

단순한 관찰이나 측정 등의 수단을 통해 현실 세계로부터 수집된 사실이나 값

의미있게 사용하기 위해서 구조화가 필요함

* 구조화: 체계적으로 조직하는 것

DBMS



* 이미지 출처: [생활코딩 MySQL 수업](#)

데이터베이스란

엑셀

DBMS

파일

데이터베이스

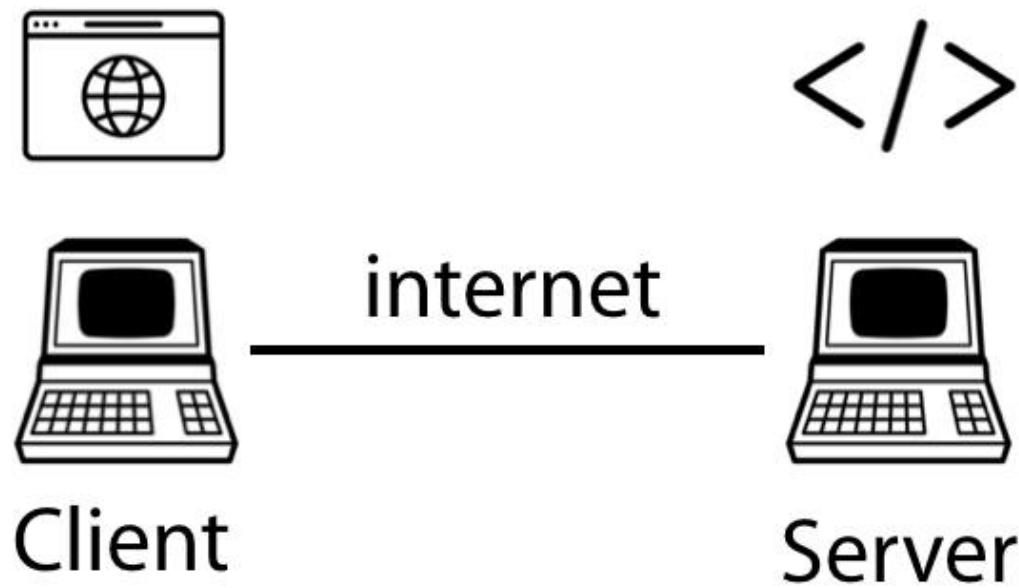
시트

테이블

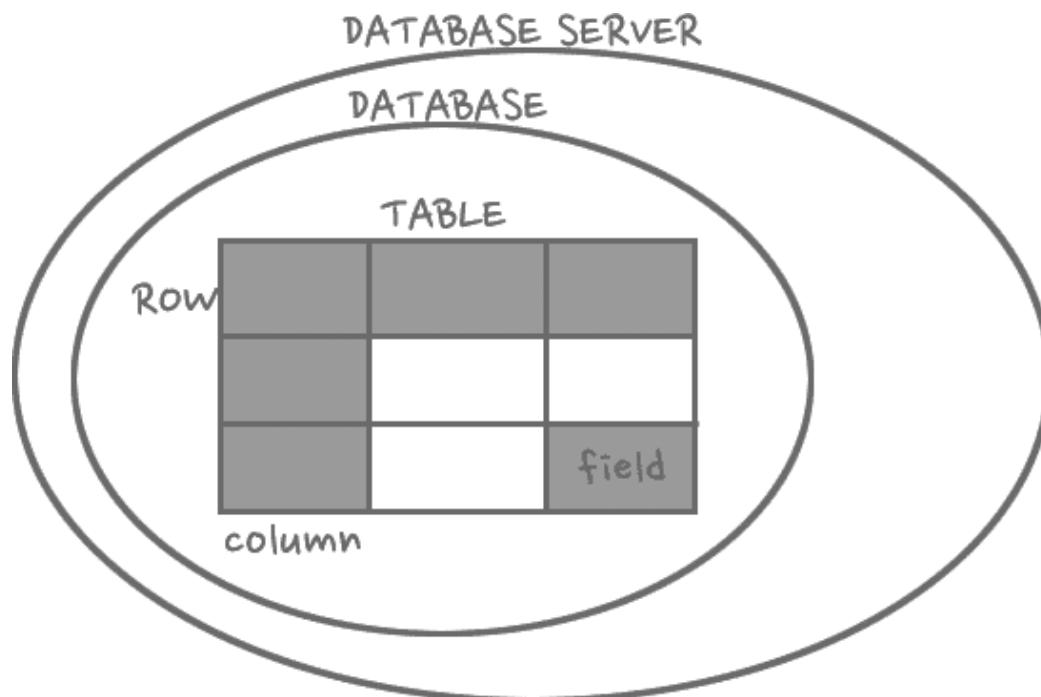
데이터베이스 클라이언트

- 대표적인 클라이언트
 - MySQL monitor
 - PHPmyAdmin
 - Navicat
 - HeidiSQL

서버와 클라이언트



테이블



* 이미지 출처: [생활코딩 MySQL 수업](#)

data.frame

```
## Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1      5.1       3.5       1.4       0.2   setosa
## 2      4.9       3.0       1.4       0.2   setosa
## 3      4.7       3.2       1.3       0.2   setosa
## 4      4.6       3.1       1.5       0.2   setosa
## 5      5.0       3.6       1.4       0.2   setosa
## 6      5.4       3.9       1.7       0.4   setosa
```

SQL

Structured Query Language

구조적 데이터 요청 언어

DBI

```
## Loading required package: devtools
```

```
## Loading required package: DBI
```

```
## Loading required package: RSQLite
```

SQLite

SQLite

SQLite is a self-contained, high-reliability, embedded, full-featured, public-domain, SQL database engine.

파일 하나로 구성하는 작고, 무료인 sql db

DBI

```
library(DBI)
library(RSQLite)
con <- dbConnect(RSQLite::SQLite(), dbname="class2.sqlite")

dbListTables(con)

dbWriteTable(con, "mtcars", mtcars, overwrite=T)
dbListTables(con)

dbReadTable(con, "mtcars")

dbRemoveTable(con, "mtcars")
dbListTables(con)

system.time(dbWriteTable(con, "member", "./recomen/membership.csv", row.names=F))
```

데이터

데이터 공유

Leek group에서 소개하는 데이터 공유 가이드

- 원시 데이터
- 정제후 데이터
- 코드북
- 변수 작성법
- 재현성

원시 데이터

최초 획득한 당시 그대로의 데이터

- 어떤 식으로든 수정을 하지 않은 상태
- 수정을 하는 과정을 함께 기록함으로써 신뢰성 확보
- 위 두 가지가 없는 경우 상황을 상상해야 함

정제후 데이터

해들리 위컴이 설명한 tidy data의 요건에 맞게 가공하여 데이터를 쉽게 다룰 수 있게 만든 상태

- 측정하는 각 변수는 하나의 열에 있어야 함
- 측정하는 각 관찰은 하나의 행에 있어야 함
- 각 종류의 변수에 대해 각 하나의 테이블이 있어야 함
- 여러 개의 테이블이 있는 경우 테이블에 합치기 위한 기준 열을 포함해야 함

코드북

데이터셋에 대해 필요한 설명을 담은 문서

- 정제후 데이터에 대해 추가적으로 필요한 설명이나 정보(단위 등)
- 정제 과정에서 사용한 방법의 설명과 사용한 이유
- 데이터가 사용된 분석에 대한 정보

데이터 원본

확보할 당시의 원시 데이터나, 항상 최신 상태를 유지하여 신뢰할 수 있는 데이터

- 커뮤니케이션 비용 감소
- 의사결정 및 활동의 기준
- 가공된 데이터의 신뢰성 확보

데이터 훑어보기

- head: 최초 6행의 데이터를 보여줌(행갯수 조절 가능)
- tail: 마지막 6행의 데이터를 보여줌(행갯수 조절 가능)
- summary: 각 컬럼의 자료형과 숫자라면 대표값을 함께 보여줌
- str: 각 컬럼의 자료형과 초기 값을 보여줌
- length: 데이터의 길이 출력(vector)
- nrow: 행 갯수 출력(data.frame)
- is.na: NA 인지 확인
- complete.cases: 값이 모두 있는지 행단위로 검사
- tibble: 최근 기법으로 재구성된 data.frame

추천 데이터 - chennel

- cusID : 5자리 숫자조합으로 구성된 고객ID
- chennel: 접속 체널
- useCnt : 사용횟수(건)

```
##      cusID       chennel        useCnt
## Min.   :    7  Length:8824      Min.   : 1.00
## 1st Qu.: 6107  Class :character 1st Qu.: 2.00
## Median : 9506  Mode  :character Median : 7.00
## Mean   : 9835                           Mean   :13.57
## 3rd Qu.:13812                           3rd Qu.:19.00
## Max.   :19382                           Max.   :240.00
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 8824 obs. of 3 variables:
## $ cusID : int 7 14 42 74 74 94 112 112 122 123 ...
## $ chennel: chr "A_MOBILE/APP" "A_MOBILE/APP" "B_MOBILE/APP" "A_MOBILE/APP" ...
## $ useCnt : int 4 1 23 1 30 14 16 1 27 10 ...
## - attr(*, "spec")=List of 2
##   ..- attr(*, "class")= chr "col_spec"
```

추천 데이터 - competitor

- cusID : 5자리 숫자조합으로 구성된 고객ID
- partner : 제휴사
- competitor: 경쟁사
- useDate : 이용년월(YYYYDD)

```
##      cusID          partner        competitor       useDate
##  Length:28159    Length:28159    Length:28159   Min.   :201501
##  Class :character Class :character Class :character 1st Qu.:201504
##  Mode  :character Mode  :character Mode  :character Median  :201507
##                                         Mean   :201507
##                                         3rd Qu.:201510
##                                         Max.   :201512
```

추천 데이터 - competitor

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 28159 obs. of 4 variables:  
## $ cusID    : chr "00002" "00051" "00077" "00077" ...  
## $ partner   : chr "D" "D" "D" "D" ...  
## $ competitor: chr "D02" "D01" "D02" "D02" ...  
## $ useDate   : int 201507 201504 201503 201506 201507 201508 201511 201510 201511 201508 ...  
## - attr(*, "spec")=List of 2  
##   ..- attr(*, "class")= chr "col_spec"
```

추천 데이터 - customer

- cusID: 5자리 숫자조합으로 구성된 고객ID
- sex : 성별
- age : 연령
 - 19세이하, 20세~24세, 25세~29세, 30세~34세, 35세~39세, 40세~44세, 45세~49세, 50세~54세, 55세~59세, 60세이상
- area : 거주지역

추천 데이터 - customer

```
##      cusID             sex            age
##  Length:19383    Length:19383    Length:19383
##  Class :character  Class :character  Class :character
##  Mode  :character  Mode  :character  Mode  :character
##      area
##  Length:19383
##  Class :character
##  Mode  :character

## Classes 'tbl_df', 'tbl' and 'data.frame': 19383 obs. of 4 variables:
## $ cusID: chr "00001" "00002" "00003" "00004" ...
## $ sex  : chr "M" "M" "M" "F" ...
## $ age  : chr "60세이상" "60세이상" "60세이상" "60세이상" ...
## $ area : chr "060" "100" "033" "016" ...
## - attr(*, "spec")=List of 2
##   ..- attr(*, "class")= chr "col_spec"
```

추천 데이터 - item

- partner : 재휴사
- cate_1 : 대분류
- cate_2 : 중분류
- cate_3 : 소분류
- cate_2_name: 중분류명
- cate_3_name: 소분류명

추천 데이터 - item

```
##   partner          cate_1          cate_2
## Length:4386    Length:4386    Length:4386
## Class :character  Class :character  Class :character
## Mode  :character  Mode  :character  Mode  :character
##   cate_3          cate_2_name     cate_3_name
## Length:4386    Length:4386    Length:4386
## Class :character  Class :character  Class :character
## Mode  :character  Mode  :character  Mode  :character
```

추천 데이터 - item

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 4386 obs. of 6 variables:  
## $ partner : chr "A" "A" "A" "A" ...  
## $ cate_1   : chr "01" "01" "01" "01" ...  
## $ cate_2   : chr "0101" "0101" "0101" "0101" ...  
## $ cate_3   : chr "A010101" "A010102" "A010103" "A010104" ...  
## $ cate_2_name: chr "일용잡화" "일용잡화" "일용잡화" "일용잡화" ...  
## $ cate_3_name: chr "위생세제" "휴지류" "뷰티상품" "일용잡화" ...  
## - attr(*, "spec")=List of 2  
##   ..- attr(*, "class")= chr "col_spec"
```

추천 데이터 - membership

- cusID : 5자리 숫자조합으로 구성된 고객ID
- memberShip: 멤버십명
- regDate : 가입년월

```

##      cusID       memberShip        regDate
##  Length:7456   Length:7456   Min.    :201210
##  Class :character  Class :character  1st Qu.:201311
##  Mode  :character  Mode  :character  Median  :201407
##                                         Mean   :201412
##                                         3rd Qu.:201504
##                                         Max.   :201512

## Classes 'tbl_df', 'tbl' and 'data.frame':    7456 obs. of  3 variables:
## $ cusID      : chr  "00011" "00021" "00037" "00043" ...
## $ memberShip: chr  "하이마트" "하이마트" "하이마트" "하이마트" ...
## $ regDate    : int  201512 201506 201306 201403 201411 201312 201506 201404 201406 201311 ...
## - attr(*, "spec")=List of 2
##   ..- attr(*, "class")= chr "col_spec"

```

추천 데이터 - tran

```

##   partner          receiptNum        cate_1        cate_2
## Length:28593030    Min. :     1    Min. : 1.00    Min. : 101
## Class :character  1st Qu.:3922474  1st Qu.: 4.00  1st Qu.: 401
## Mode  :character  Median :7167787  Median :11.00  Median :1102
##                           Mean   :6447881  Mean   :18.37  Mean   :1840
##                           3rd Qu.:9116336  3rd Qu.:18.00  3rd Qu.:1808
##                           Max.  :11096601  Max.  :92.00  Max.  :9206
##   cate_3            cusID       storeCode        date
## Length:28593030    Min. :     1    Min. : 1.00  Min. :20140101
## Class :character  1st Qu.: 5206  1st Qu.: 16.00 1st Qu.:20140711
## Mode  :character  Median :10104  Median : 44.00  Median :20150110
##                           Mean   : 9904  Mean   : 92.26  Mean   :20145817
##                           3rd Qu.:14638  3rd Qu.:110.00 3rd Qu.:20150703
##                           Max.  :19383  Max.  :593.00  Max.  :20151231
##   time             amount
## Min.  : 0.00  Min.  :      1
## 1st Qu.:14.00 1st Qu.:    2050
## Median :17.00 Median :    4290
## Mean   :16.71 Mean   :   23678
## 3rd Qu.:19.00 3rd Qu.:   9900
## Max.  :23.00  Max.  :101330000

```

추천 데이터 - tran

```
## Classes 'tbl_df', 'tbl' and 'data.frame': 28593030 obs. of 10 variables:  
## $ partner : chr "B" "B" "B" "B" ...  
## $ receiptNum: int 8664000 8664000 8664000 8664000 8664001 8664001 8664002 8664002 8664002 8664003 ...  
## $ cate_1   : int 15 16 16 18 5 15 10 43 54 5 ...  
## $ cate_2   : int 1504 1601 1602 1803 509 1501 1003 4301 5403 504 ...  
## $ cate_3   : chr "B150401" "B160101" "B160201" "B180301" ...  
## $ cusID    : int 17218 17218 17218 17218 17674 17674 14388 14388 14388 15773 ...  
## $ storeCode: int 44 44 44 44 44 44 44 44 44 44 ...  
## $ date     : int 20140222 20140222 20140222 20140222 20140222 20140222 20140222 20140222 20140222 20140222 ...  
## $ time     : int 20 20 20 20 22 22 23 23 23 21 ...  
## $ amount   : int 2420 1070 8060 6000 1120 1200 5290 5960 9900 970 ...  
## - attr(*, "spec")=List of 2  
##   ..- attr(*, "class")= chr "col_spec"
```

클라우드 서비스

클라우드 서비스 소개

클라우딩 컴퓨팅은 사용자의 환경 밖에서 서비스로서 제공된 확장 가능한 컴퓨팅 자원을 사용한 양에 따라 비용을 지불하고 사용하는 것

출처: P. Changanti, 가상 인프라용 클라우드 서비스, Part 1: IaaS(Infrastructure as a Service) 및 Eucalyptus

구글 클라우드

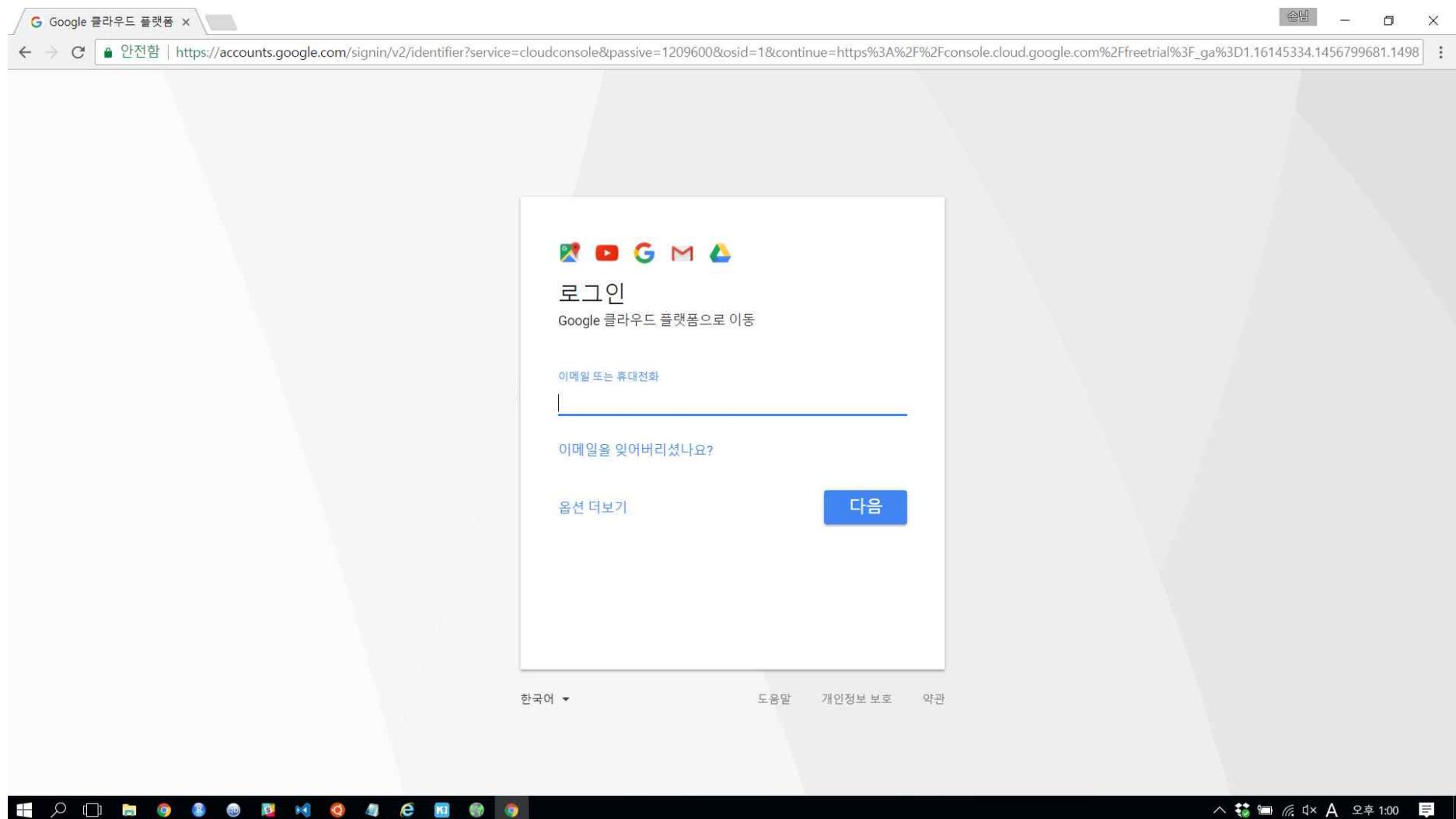


Google Cloud Platform

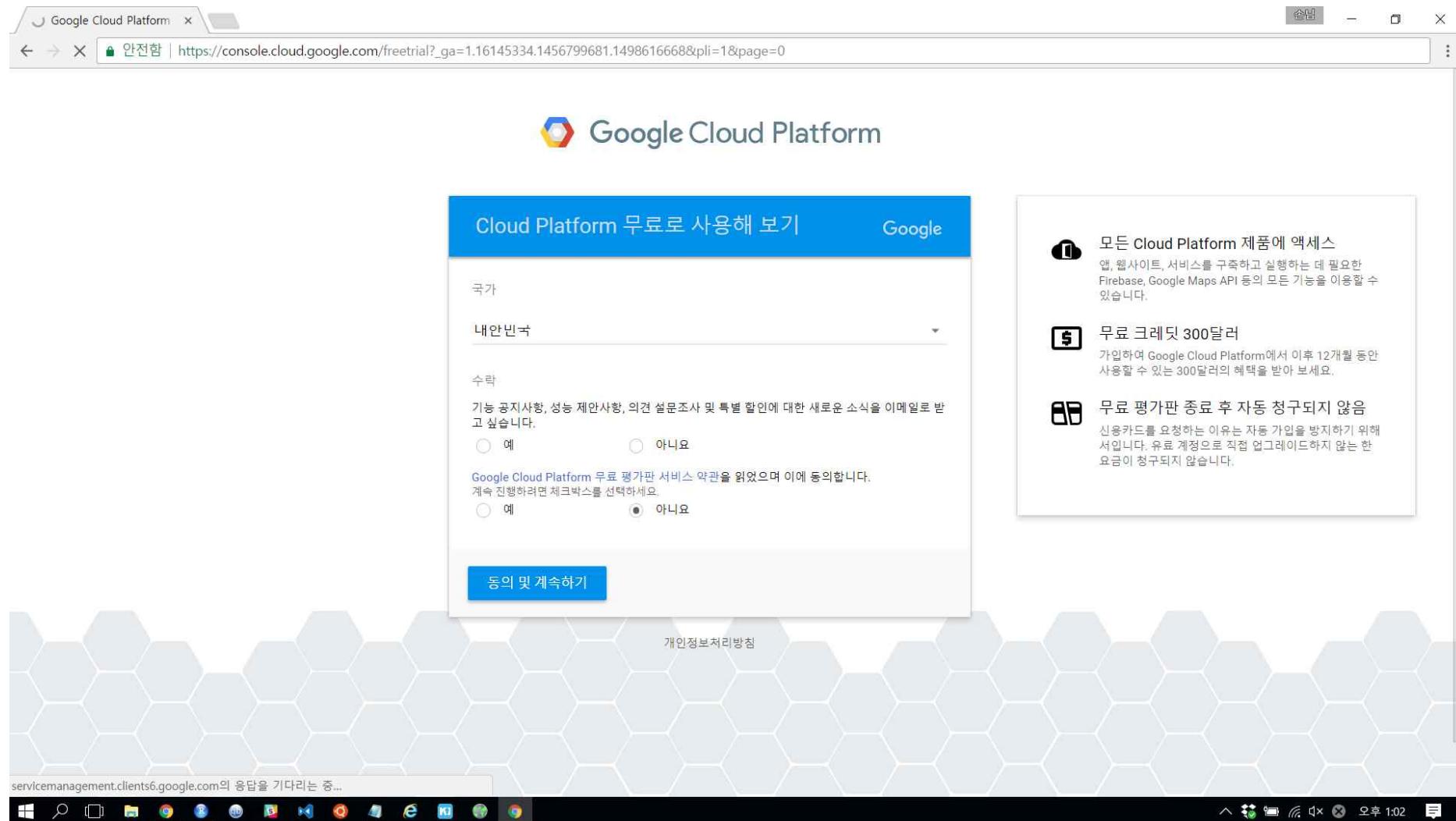
구글 클라우드 시작하기

<https://cloud.google.com/>

구글 클라우드 로그인

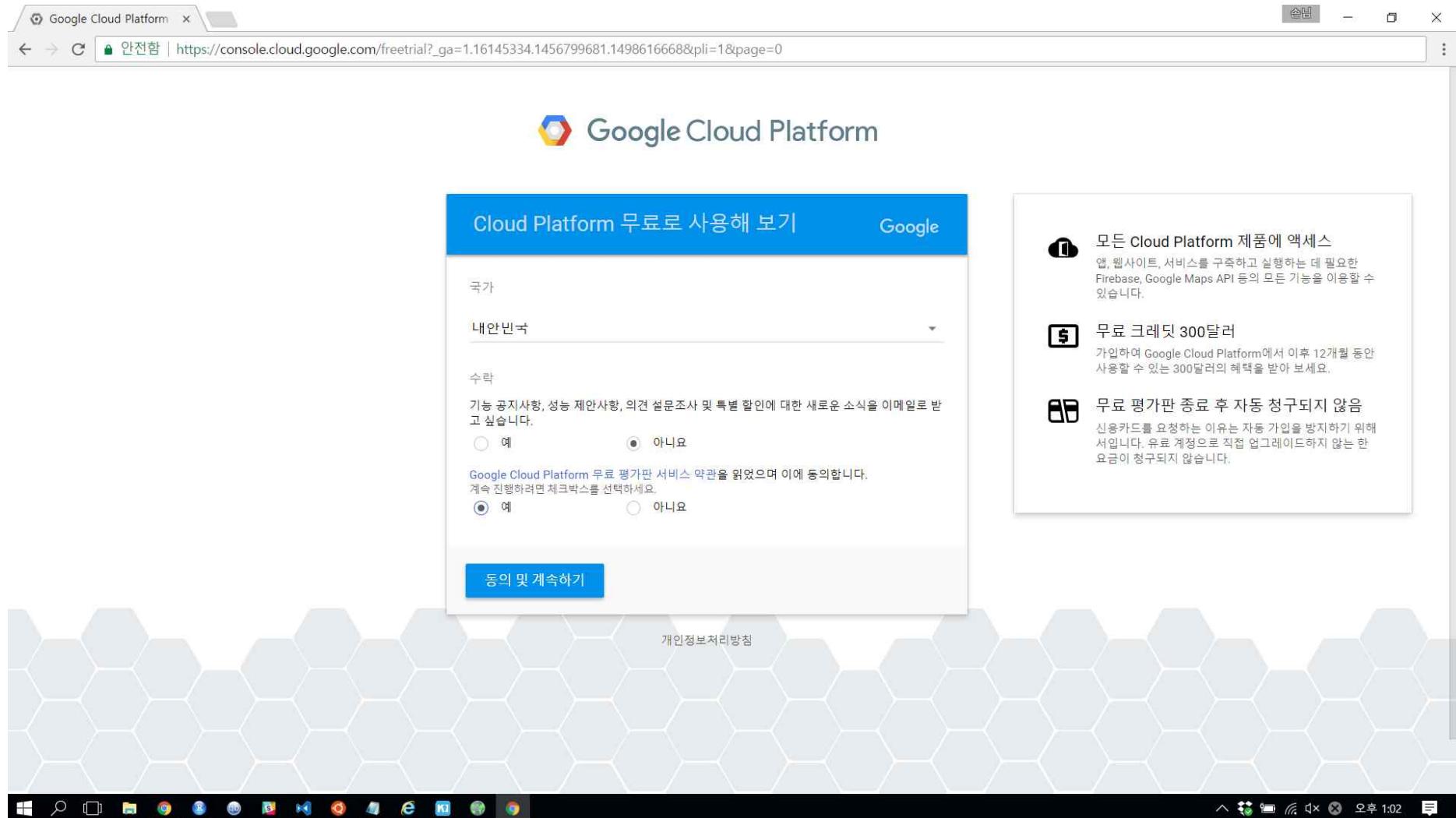


구글 클라우드 설문

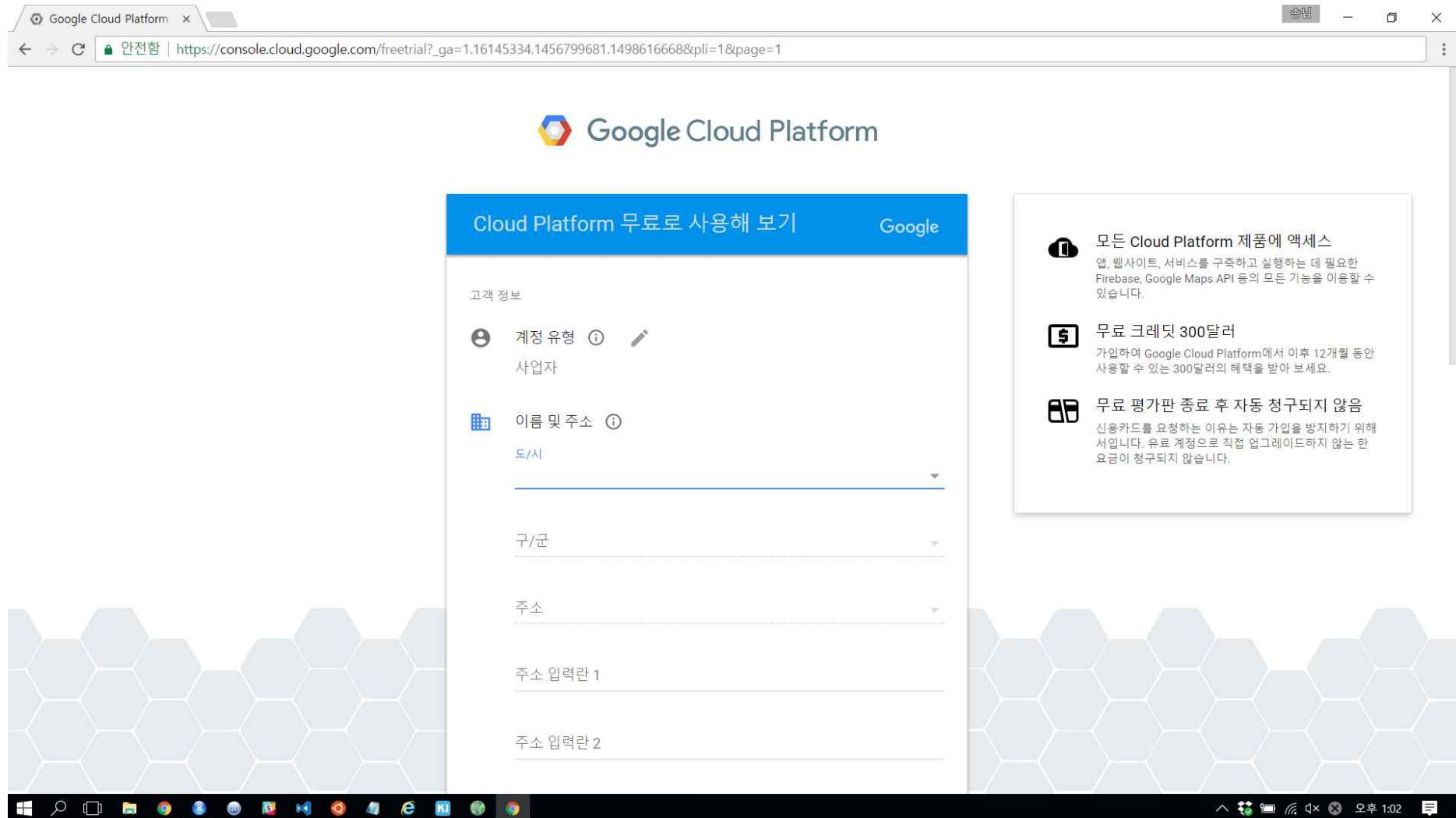


39/92

구글 클라우드 설문 선택



구글 클라우드 카드 등록



구글 클라우드 프로젝트 생성

The screenshot shows the Google Cloud Platform dashboard for the project 'konlper'. The left sidebar lists various services: Compute Engine, Datastore, Storage, SQL (selected), Spanner, StackDriver, Monitoring, Debug, Metrics, Logs, and Container Registry. The main content area has tabs for '대시보드' (Dashboard) and '활동' (Activities). The '대시보드' tab is active, displaying sections for '프로젝트 정보', 'Compute Engine', 'SQL', and '추적' (Metrics). The 'Compute Engine' section shows a CPU usage chart from June 28, 2017, at 12:30 to 1:06, with a value of 3.083. The 'SQL' section shows storage usage: 1G, 768M, 512M, and 256M. The '추적' section indicates no data for the past 7 days. On the right, there are sections for 'Google Cloud Platform 상태' (Status), '결제' (Billing), '오류 보고' (Error Reporting), and '뉴스' (News). The status section shows all services are up. The billing section shows a balance of \$0.00. The error reporting section notes no errors found. The news section links to information about enterprise identity.

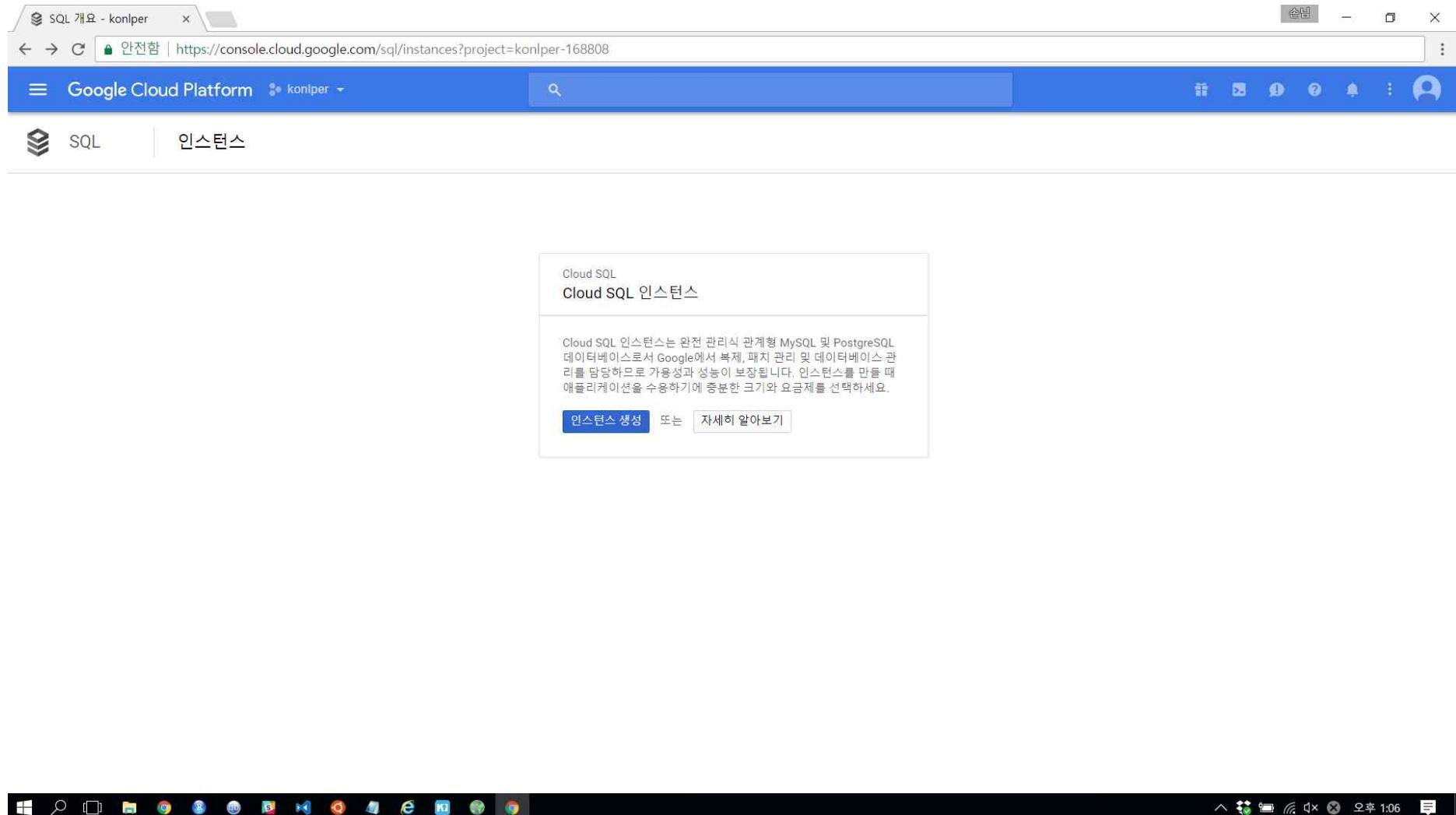
42/92

R MySQL 연결

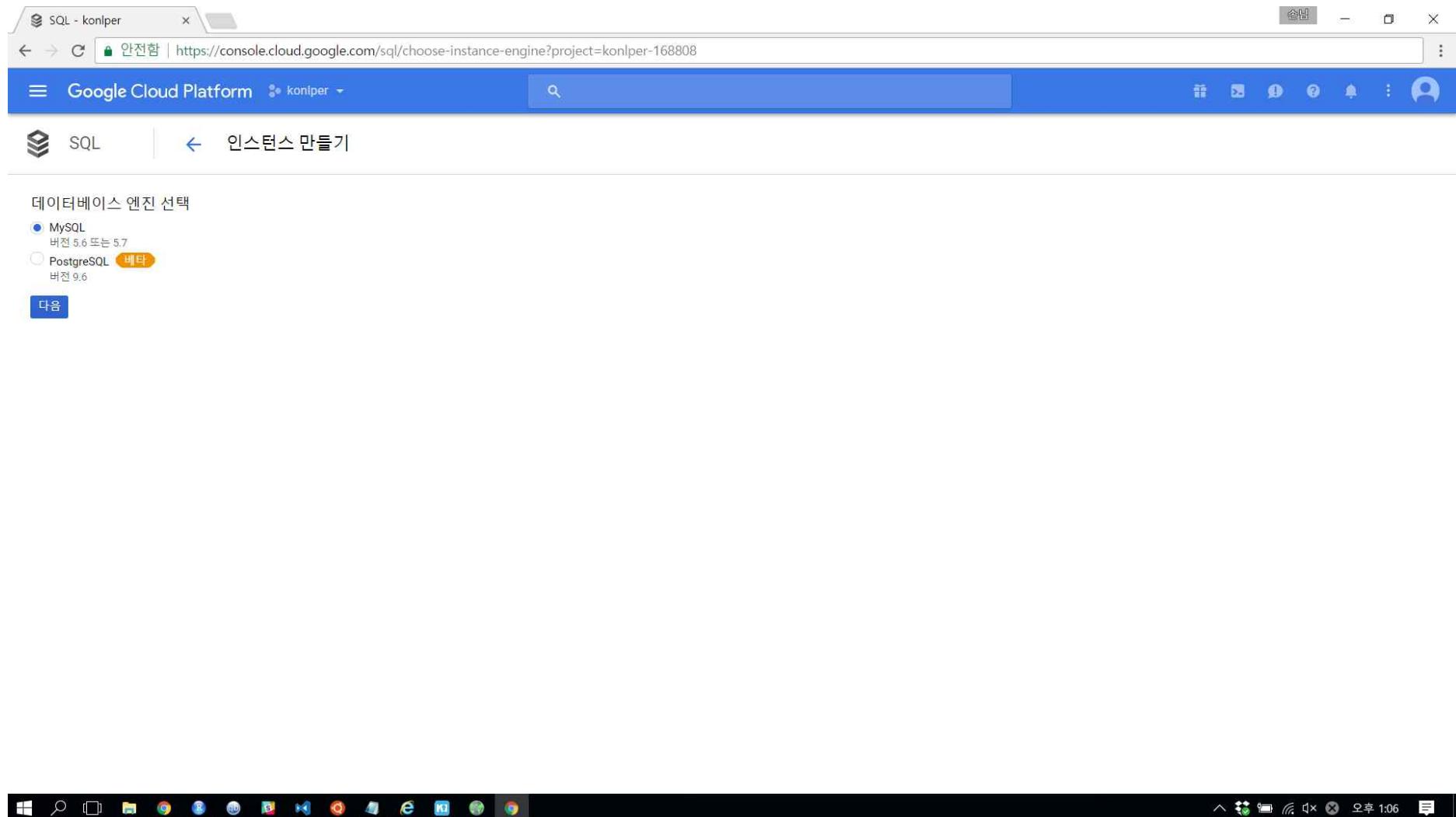
DBI로 MySQL을 연결하는 패키지와 사용

```
library(RMySQL)
con <- dbConnect(MySQL(),
                  user = user,
                  password = pw,
                  host = host,
                  dbname = "recom")
dbListTables(conn = con)
dbWriteTable(conn = con, name = 'tran', value = "./recomen/tran.csv")
dbReadTable(conn = con, name = "tran")
```

MySQL 인스턴스 생성



엔진선택



MySQL 2세대 선택

The screenshot shows the Google Cloud Platform SQL interface. The URL in the browser is <https://console.cloud.google.com/sql/pre-create?project=konlper-168808>. The page title is "MySQL 인스턴스 유형 선택". On the left, there's a sidebar with "SQL" selected. The main content area has two sections: "MySQL 2세대(권장)" and "MySQL 1세대(이전)". The "MySQL 2세대(권장)" section is highlighted with a blue border. It contains a summary: "저렴한 비용으로 높은 성능과 큰 저장용량이 제공됩니다." followed by a bulleted list: ● 1세대 처리량의 최대 7배 및 저장용량 크기 20배 ● 대부분의 사용 사례에서 1세대보다 비용 낮음 ● 고가용성 장애 조치를 추가하고 복제본을 읽을 수 있는 옵션 제공 ● 백업 기간 및 유지관리 기간 구성 가능 ● MySQL 5.6 및 5.7만 지원됩니다. A blue button labeled "2세대 선택" is at the bottom of this section. The "MySQL 1세대(이전)" section is shown in a greyed-out state with a white border. It contains the text: "기본 성능과 저장용량 크기를 제공하는 이전 버전의 Cloud SQL입니다. MySQL 5.7을 지원하지 않습니다." and a greyed-out button labeled "1세대 선택".



인스턴스 설정

The screenshot shows the 'MySQL 2세대 인스턴스 만들기' (Create MySQL 2nd Gen Instance) page in the Google Cloud Platform. The page includes fields for '인스턴스 ID' (Instance ID), '데이터베이스 버전' (Database Version), '위치' (Location) with '지역' (Region) set to 'us-central1' and '영역' (Zone) set to '자동 선택' (Auto-select), '머신 유형' (Machine Type) showing 'db-n1-standard-1' with 1 vCPU and 3.75GB memory, and '네트워크 처리량(MB/초)' (Network Throughput) set to 250/2,000. The '저장소 유형' (Storage Type) section has 'SSD(권장)' (SSD (Recommended)) selected. The browser address bar shows the URL: https://console.cloud.google.com/sql/create-instance-mysql?project=konlper-168808.

접속 허용 IP 설정

The screenshot shows the Google Cloud Platform SQL interface for creating a MySQL 2nd Gen instance. The page has the following sections:

- 유지관리 기간**: Includes dropdowns for '시스템 자동 선택 기간' and '유지관리 시점'.
- 루트 비밀번호**: A field for setting the root user's password, with a '생성' (Generate) button and a checkbox for '비밀번호 없음'.
- Cloud SQL 플래그**: A section with a '+ 항목 추가' (Add item) button.
- 승인된 네트워크**: A section for whitelisting IPv4 addresses, with a '+ 네트워크 추가' (Add network) button.
- Bottom Buttons**: '생성' (Create) and '취소' (Cancel) buttons.

접속 허용 IP 설정

The screenshot shows the Google Cloud Platform SQL interface for creating a MySQL instance. The URL in the browser is <https://console.cloud.google.com/sql/create-instance-mysql?project=konlper-168808>. The main page title is "MySQL 2세대 인스턴스 만들기". A sub-section titled "Cloud SQL 플래그" contains a button "+ 항목 추가". Below this, a section for "승인된 네트워크" asks to add an IPv4 address to the instance's network. A modal window titled "새 네트워크" (New Network) is open, asking for a name ("이름 (선택사항)" with "없음" selected) and a CIDR range ("네트워크" with "예: 199.27.25.0/24"). At the bottom of the modal are "완료" and "취소" buttons, and a "네트워크 추가" button. The status bar at the bottom shows the URL en.wikipedia.org/wiki/Classless_Inter-Domain_Routing#CIDR_notation and the system time "오후 1:07".

49/92

현재 IP 확인하기

A screenshot of a Google search results page for the query "myip". The search bar shows "myip". The results include:

- What Is My IP Address - IP Address Tools and Info - WhatIsMyIP.com ®**
https://www.whatismyip.com/ ▾ 이 페이지 번역하기
We provide IP address tools that allow users to perform an Internet Speed Test, IP address lookup, proxy detection, IP Whois Lookup, and more.
- My IP Information**
My IP Information tool shows your ip address, city, state, country ...
whatismyip.com 검색결과 더보기 »
- IP Address Lookup**
The IP Address Lookup tool includes the following IP ...
- What Is My IP Address? IP Address Tools and More**
whatismyipaddress.com/ ▾ 이 페이지 번역하기
IP address lookup, location, proxy detection, email tracing, IP hiding tips, blacklist check, speed test, and forums. Find, get, and show my IP address.
- What is my IP address? - IP Location**
https://www.iplocation.net/find-ip-address ▾ 이 페이지 번역하기
This webpage displays public IP address of your computer or router assigned by your ISP.
- What's My IP Address? | Online Privacy and Security Tool - ExpressVPN**
https://www.expressvpn.com/what-is-my-ip ▾ 이 페이지 번역하기
IP address lookup. What do others know about your location? See the true IP address of your VPN or proxy server. Learn to hide your IP address in 2 minutes.
- What's My IP Address? Networking Tools & More**
www.whatismyip.org/ ▾ 이 페이지 번역하기
Your IP Address plus Port Scanners, Traceroute, HTTP Compression Test, Ping, Whois, DNS, IP Geo Location, Password Generator and many more tools and ...

whatismyip

<https://www.whatismyip.com/>

ip 입력

The screenshot shows a browser window for the Google Cloud Platform SQL service. The URL is <https://console.cloud.google.com/sql/create-instance-mysql?project=konlper-168808>. The main page displays a step to 'MySQL 2세대 인스턴스 만들기' (Create MySQL 2nd Gen Instance). A sub-modal window titled '새 네트워크' (New Network) is open, asking for an IP address. The input field contains '112.217.214.251'. The bottom of the screen shows a Windows taskbar with various icons.

52/92

root 계정 비밀번호 생성

The screenshot shows the Google Cloud Platform SQL interface for creating a MySQL instance. The page title is "MySQL 2세대 인스턴스 만들기". The "Root 비밀번호" section contains a password input field and a "생성" (Create) button. The "Cloud SQL 플래그" section has a "항목 추가" (Add item) button. The "승인된 네트워크" section lists "112.217.214.251" with a "제거/지우기" (Delete) link. At the bottom are "생성" and "취소" buttons, and a Windows taskbar at the bottom.

SQL 인스턴스 생성중

The screenshot shows a browser window with the URL <https://console.cloud.google.com/sql/instances?project=konlper-168808>. The page is titled "Google Cloud Platform" and "konlper". The main content area is titled "SQL" and "인스턴스". It displays a table with one row:

인스턴스 ID	유형	IP 주소	고가용성	위치	사용된 저장용량
fctestp	MySQL 2세대	-	-	us-central1	-

The status of the instance is listed as "생성 중" (Creating) in the "상태" column. The bottom of the screen shows a Windows taskbar with various icons and the system tray.

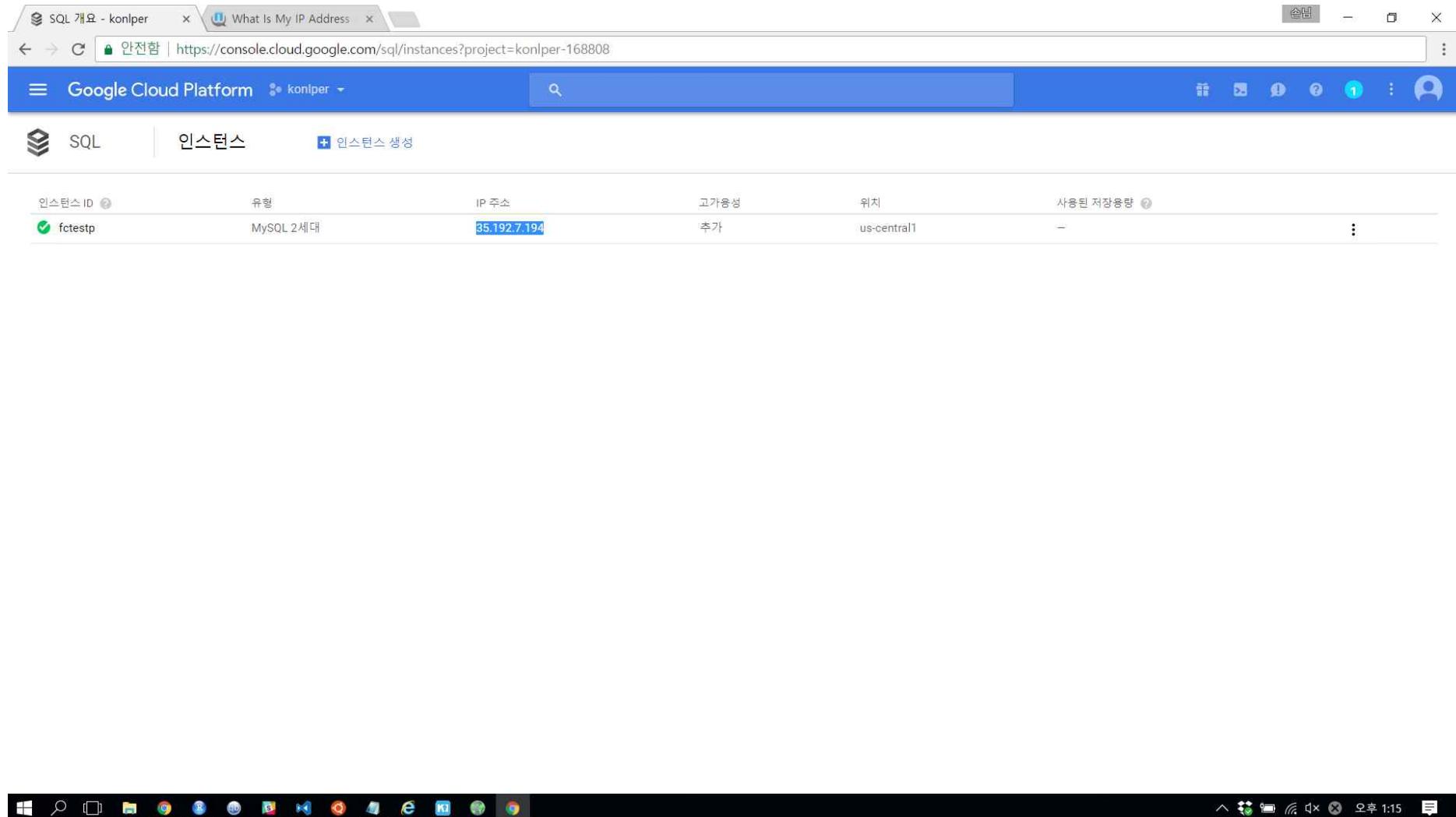
RMySQL로 연결하기

- host 주소 파악

The screenshot shows a Windows desktop environment. At the top, there is a taskbar with several icons: Start, Search, Task View, File Explorer, Google Chrome, FileZilla, and others. A system tray icon for battery status is also visible. The main window is a web browser displaying the Google Cloud Platform SQL Instances page. The URL in the address bar is <https://console.cloud.google.com/sql/instances?project=konlper-168808>. The page title is "SQL 개요 - konlper". The content area shows a table of instances:

인스턴스 ID	유형	IP 주소	고가용성	위치	사용된 저장용량
fctestp	MySQL 2세대	35.192.7.194	추가	us-central1	-

RMySQL로 연결하기



The screenshot shows the Google Cloud Platform SQL Instances page. The browser tabs include 'SQL 개요 - konlper' and 'What Is My IP Address'. The address bar shows the URL <https://console.cloud.google.com/sql/instances?project=konlper-168808>. The main content area displays a table with one row of data:

인스턴스 ID	유형	IP 주소	고가용성	위치	사용된 저장 용량
fctestp	MySQL 2세대	35.192.7.194	추가	us-central1	-

The Windows taskbar at the bottom shows various pinned icons and the system tray.

데이터베이스 만들기



데이터베이스 만들기

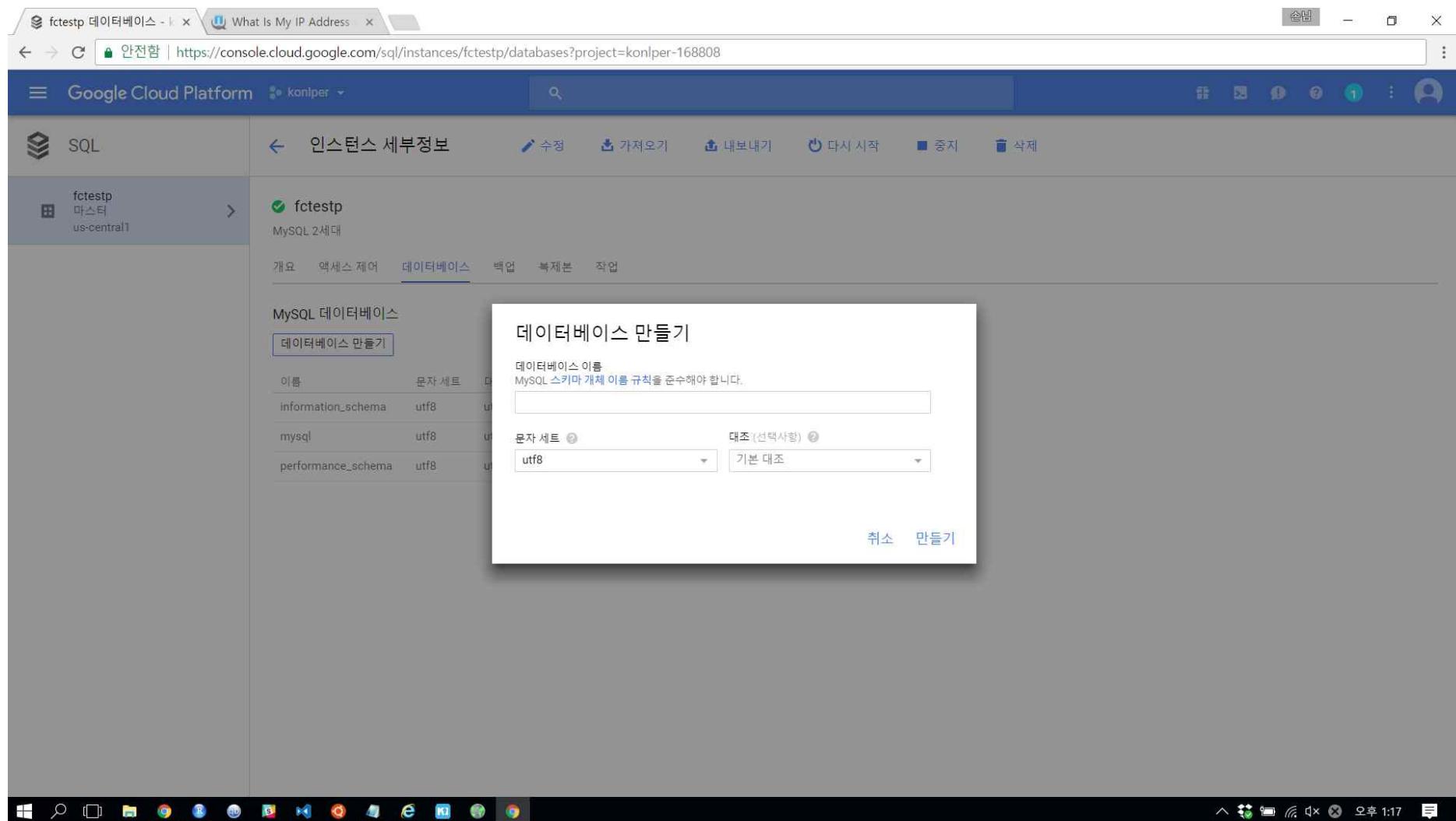
The screenshot shows the Google Cloud Platform SQL Instances page. On the left, a sidebar lists 'SQL' and an instance named 'fctestp' under the '마스터' section of the 'us-central1' region. The main content area displays the 'fctestp' instance details. At the top, there are buttons for '수정' (Edit), ' 가져오기' (Import), '내보내기' (Export), '다시 시작' (Restart), '중지' (Stop), and '삭제' (Delete). Below these are tabs for '개요' (Overview), '액세스 제어' (Access Control), '데이터베이스' (Database), '백업' (Backup), '복제본' (Replica), and '작업' (Operations). A dropdown menu for '저장소 사용량' (Storage Usage) is open, showing options like 1G, 768M, 512M, and 256M. A timeline at the bottom shows storage usage from 6월 28일 오후 12:30 to 6월 28일 오후 1:00, with a total usage of 1.14G. Below the timeline, there's a '로그 기록' (Log History) section with a button for '오류 로그 보기' (View Error Log) and a '이 인스턴스에 연결' (Connect to this instance) section with buttons for 'Cloud Shell을 사용해 연결' (Connect via Cloud Shell) and '모든 연결 방법 보기' (View all connection methods).

58/92

데이터베이스 만들기

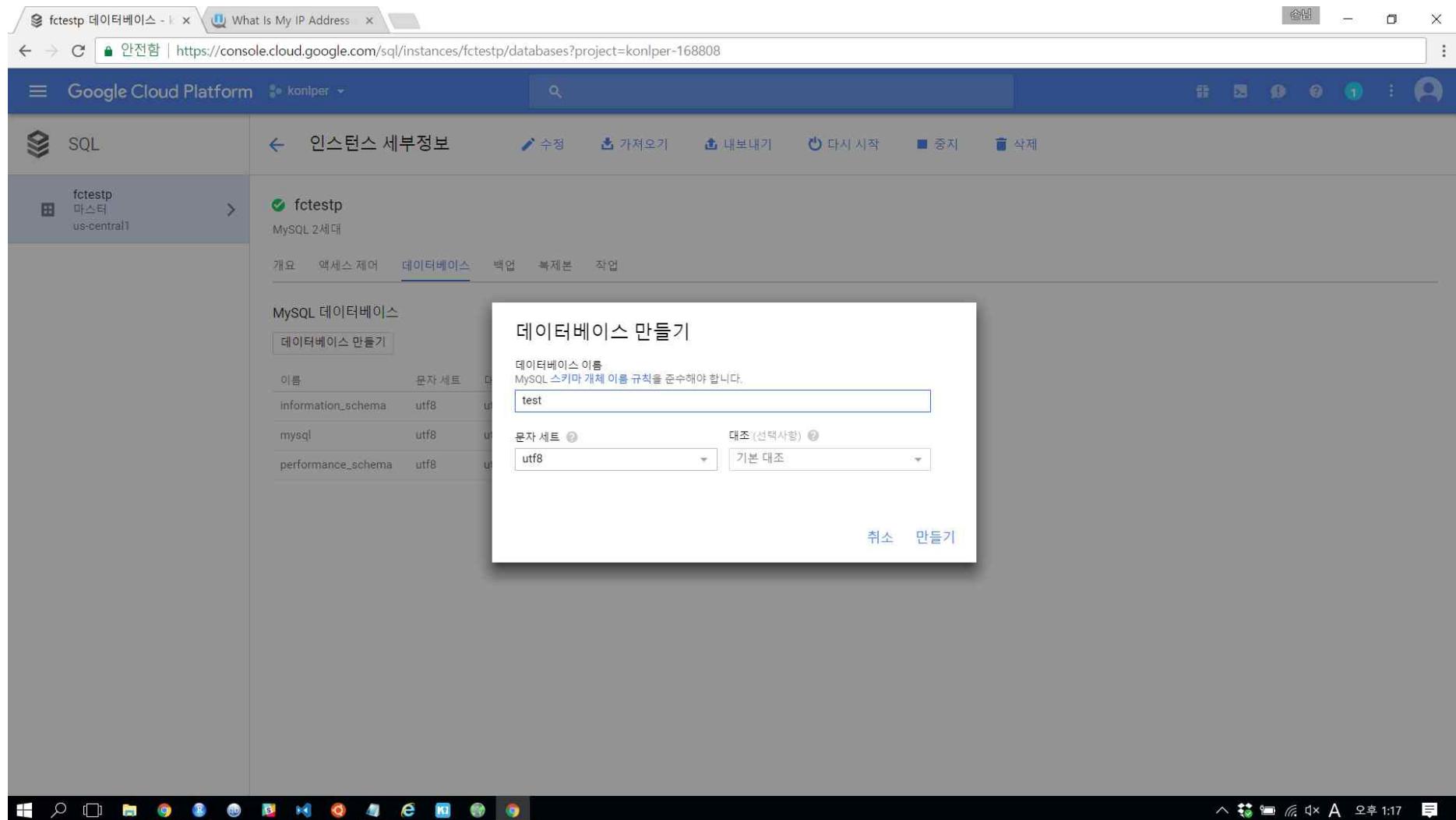
The screenshot shows the Google Cloud Platform SQL interface. On the left, a sidebar lists 'fctestp' instances under the '마스터' section. The main panel displays the 'Instances 세부정보' (Instance Details) for 'fctestp'. At the top, there are buttons for '수정' (Edit), '가져오기' (Import), '내보내기' (Export), '다시 시작' (Restart), '중지' (Stop), and '삭제' (Delete). Below these are tabs for '개요' (Overview), '액세스 제어' (Access Control), '데이터베이스' (Database), '백업' (Backup), '복제본' (Replica), and '작업' (Jobs). The '데이터베이스' tab is selected. A sub-section titled 'MySQL 데이터베이스' shows a table with three rows: 'information_schema', 'mysql', and 'performance_schema'. A button labeled '데이터베이스 만들기' (Create Database) is visible above the table. The bottom of the screen shows a Windows taskbar with various icons.

데이터베이스 만들기

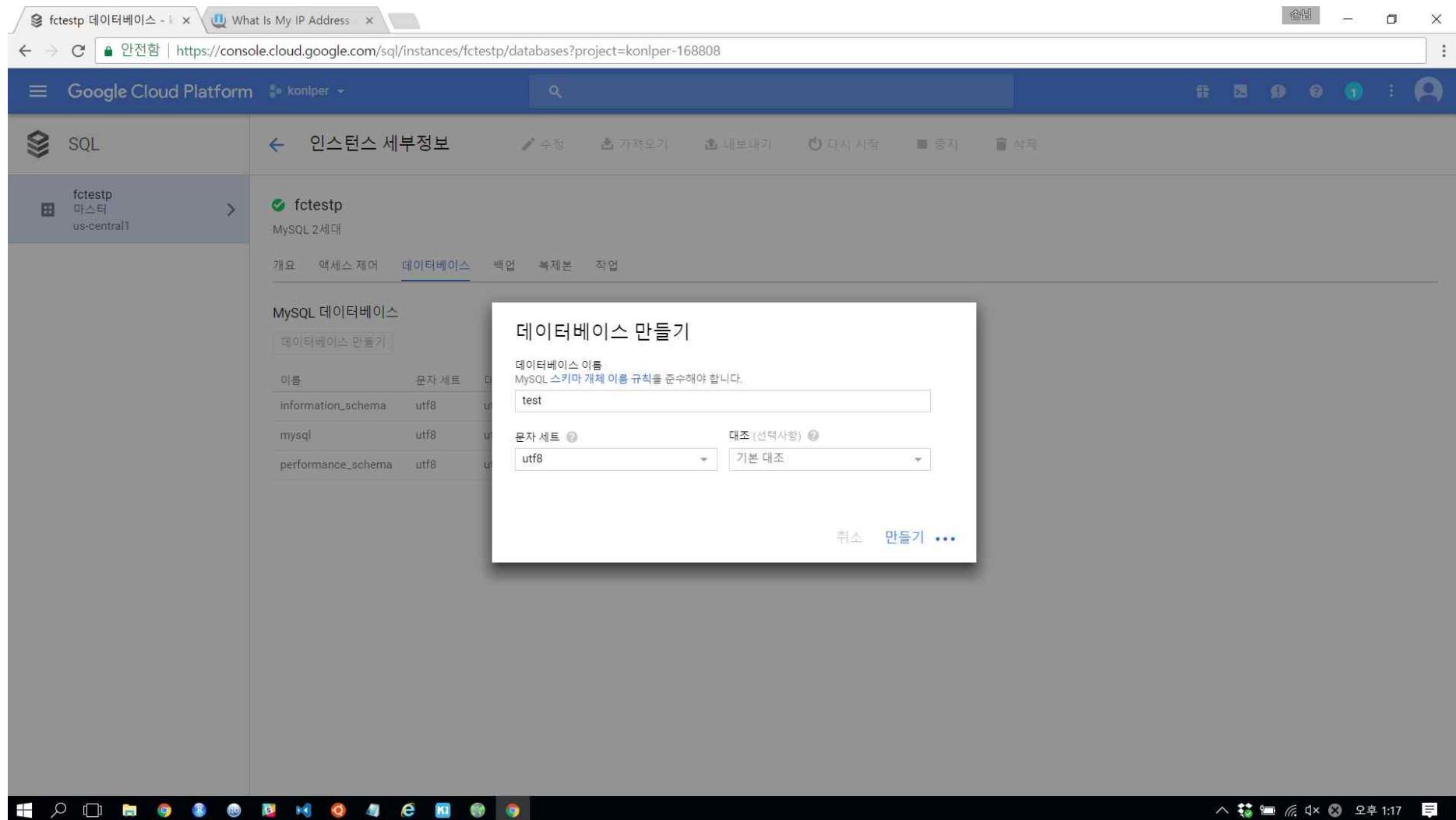


60/92

데이터베이스 만들기



데이터베이스 만들기



host와 dbname 설정

The screenshot shows the Google Cloud Platform SQL interface. On the left, a sidebar lists 'fctestp' instances under '마스터 us-central1'. The main panel displays the 'fctestp' instance details, with the '데이터베이스' tab selected. It shows a table of existing MySQL databases:

이름	문자 세트	대조	유형
information_schema	utf8	utf8_general_ci	시스템
mysql	utf8	utf8_general_ci	시스템
performance_schema	utf8	utf8_general_ci	시스템
test	utf8	utf8_general_ci	사용자

A success message '데이터베이스를 생성했습니다.' is displayed at the bottom left. The taskbar at the bottom shows various application icons.

실습 진행

The screenshot shows the RStudio interface with the following components:

- Code Editor:** Displays a script named `codeForClass2.R` containing R code for database operations.
- Console:** Shows the R version information and the start of the R command-line interface.
- Environment:** Shows the current project structure and files.
- Files:** Shows the contents of the project directory, including files like `.gitignore`, `codeForClass2.R`, and `dabrp_classnote2.Rproj`.

```

~/project/dabrp_classnote2 - master - RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Run Source
.gitignore x class2.Rmd x codeForClass2.R x
1 if (!require(devtools)) install.packages("devtools")
2 if (!require(DBI)) devtools::install_github("rstats-db/DBI")
3 if (!require(RSQLite)) devtools::install_github("rstats-db/RSQLite")
4 if (!require(RMySQL)) devtools::install_github("rstats-db/RMySQL")
5
6
7 # user<- "root"
8 # pw<- "XXXXXXXXXXXXXX"
9 # host<- 'XXX.XXX.XXX.XXX'
10
11 # save(user,pw,host,file = "./gsql.RData")
12
13 load("./gsql.RData")
14
15 library(RMySQL)
16 con <- dbConnect(MySQL(),
17                   user = user,
18                   password = pw,
19                   host = host,
20                   dbname = "test")
21
22 dbListTables(conn = con)
23 dbWriteTable(conn = con, name = 'Test', value = as.data.frame(iris))
24 dbReadTable(conn = con, name = "Test")
25
10:26 [Top Level] R Script
Console ~/project/dabrp_classnote2/
R version 3.4.0 (2017-04-21) -- "You Stupid Darkness"
Copyright (C) 2017 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> |
```

R bigquery 연결

query_exec가 첫 실행이면 브라우저에서 권한 확인을 합니다.

```
library(bigrquery)
project <- "konlper-168808"
sql <- "SELECT * FROM [konlper-168808:recom.chennel] LIMIT 5"
query_exec(sql, project = project)
```

bigquery 인스턴스 생성

The screenshot shows the Google Cloud Platform SQL Instances page. The URL in the address bar is <https://console.cloud.google.com/sql/instances/fctestp/databases?project=konlper-168808>. The main content area displays the details of the MySQL instance 'fctestp'. The 'MySQL 데이터베이스' section lists several databases:

이름	문자 세트	데조	유형
information_schema	utf8	utf8_general_ci	시스템
mysql	utf8	utf8_general_ci	시스템
performance_schema	utf8	utf8_general_ci	시스템
recom	utf8	utf8_general_ci	사용자
test	utf8	utf8_general_ci	사용자

The left sidebar shows the navigation menu with 'BigQuery' selected. The bottom status bar shows the URL <https://console.cloud.google.com/bigquery?project=konlper-168808>.

66/92

bigquery 첫 화면

The screenshot shows the Google BigQuery web interface. At the top, there's a navigation bar with tabs for 'fctestp 데이터베이스' and 'Google BigQuery'. The URL in the address bar is <https://bigquery.cloud.google.com/welcome/konper-168808>. The main content area is titled 'Welcome to BigQuery!' and contains a brief introduction about the service's capabilities. Below this, there's a section for 'To get started, try one of the following options:' with a list of links. On the left side, there's a sidebar with a 'COMPOSE QUERY' button, 'Query History', 'Job History', a search bar for 'Filter by ID or label' (with 'konper' selected), and a list of 'Public Datasets' including 'bigquery-public-data:hacker_news', 'bigquery-public-data:noaa_gsod', etc. The bottom of the screen shows a Windows taskbar with various icons and the system tray.

dataset 만들기

The screenshot shows the Google BigQuery web interface. At the top, there's a navigation bar with tabs for 'fcstestp 데이터베이스' and 'Google BigQuery'. The URL in the address bar is <https://bigquery.cloud.google.com/welcome/konper-168808>. On the left, there's a sidebar with 'COMPOSE QUERY' button, 'Query History', 'Job History', and a search/filter box for 'konper'. Below that, it says 'No datasets found in this project.' and 'Please create a dataset or select a new project from the menu above.' A dropdown menu is open over the 'konper' project name, listing options: 'Create new dataset', 'Switch to project', and 'Refresh'. The main content area is titled 'Welcome to BigQuery!' and explains that it's a web service for interactive analysis of massive datasets. It suggests reading the 'BigQuery Quickstart guide' or running a query against sample data by clicking 'Compose Query'. It also mentions creating a dataset by dragging data into a table using the 'Import' menu. A section for 'Public Datasets' lists several datasets like 'bigquery-public-data:hacker_news', 'bigquery-public-data:noaa_gsod', etc.

68/92

The screenshot shows the Google BigQuery web interface. On the left, there's a sidebar with 'COMPOSE QUERY' button, 'Query History', and 'Job History'. Below that is a search/filter bar with 'Filter by ID or label' and a dropdown set to 'konlper'. A message says 'No datasets found in this project. Please create a dataset or select a new project from the menu above.' Under 'Public Datasets', a list includes 'bigquery-public-data:hacker_news', 'bigquery-public-data:noaa_gsod', 'bigquery-public-data:samples', 'bigquery-public-data:usa_names', 'gdeilt-bq:hathitrustbooks', 'gdeilt-bq:internetarchivebooks', 'lookerdata:cdc', 'nyc-tlc:green', and 'nyc-tlc:yellow'. The main area has a 'Welcome to BigQuery!' message and a list of options to get started. A 'Create Dataset' dialog box is open in the center, prompting for 'Dataset ID' (set to 'recom'), 'Data location' (set to '(unspecified)'), and 'Data expiration' (set to 'In 60 days'). At the bottom of the dialog are 'OK' and 'Cancel' buttons. The browser's address bar shows 'fcptest 데이터베이스 - Google BigQuery' and the URL 'https://bigquery.cloud.google.com/welcome/konlper-168808'. The system tray at the bottom shows various icons and the time '오후 1:23'.

table 생성

The screenshot shows the 'Create Table' interface in Google BigQuery. The 'Source Data' section is set to 'Create from source' with 'File upload' selected for location and 'CSV' for file format. The 'Destination Table' section shows 'recom' as the table name and 'Native table' as the type. In the 'Schema' section, there is one field defined with the name field, type STRING, and mode NULLABLE. The 'Options' section includes settings for field delimiter (Comma), header rows to skip (0), number of errors allowed (0), and various row-level validation options like quoted newlines and jagged rows.

Google BigQuery

Compose Query

Query History

Job History

Filter by ID or label

konper

recom

Public Datasets

- bigquery-public-data:hacker_news
- bigquery-public-data:noaa_gsod
- bigquery-public-data:samples
- bigquery-public-data:usa_names
- gdelt-bq:hathitrustbooks
- gdelt-bq:internetarchivebooks
- lookerdata:cdc
- nyc-tlc:green
- nyc-tlc:yellow

Create Table

Source Data Create from source Create empty table

Repeat job Select Previous Job

Location File upload Choose file No file chosen

File format CSV

Destination Table

Table name recom Destination table name

Table type Native table

Schema Automatically detect

Name	Type	Mode
	STRING	NULLABLE

Add Field Edit as Text

Options

Field delimiter Comma Tab Pipe Other

Header rows to skip 0

Number of errors allowed 0

Allow quoted newlines

Allow jagged rows

Ignore unknown values

큰 데이터는 storage를 통해 업로드

The screenshot shows the 'Create Table' page in the Google BigQuery web UI. The left sidebar shows 'Query History' and 'Job History' under the project 'konlper'. A dropdown menu for 'konlper' is open, showing a single item 'recom'. Below it, the 'Public Datasets' section lists various datasets like 'bigquery-public-data:hacker_news' and 'bigquery-public-data:usa_names'. The main form is titled 'Create Table' and has the following fields:

- Source Data:** 'Create from source' is selected. 'Repeat job' dropdown is set to 'Select Previous Job'. 'File upload' button is highlighted, and a file named 'tran.csv (1484069730 bytes)' is selected.
- Location:** A note states: "Uploads from the BigQuery web UI are limited to 10 MB. For larger data sizes, please load data from Google Cloud Storage."
- File format:** 'CSV' is selected.
- Destination Table:** 'Table name' is set to 'recom'. 'Table type' is 'Native table'.
- Schema:** 'Automatically detect' is checked. A table structure is shown with a single row:

Name	Type	Mode
	STRING	NULLABLE

'Edit as Text' link is present.
- Options:** 'Field delimiter' is 'Comma'. 'Header rows to skip' is '0'. 'Number of errors allowed' is '0'. 'Allow quoted newlines' and 'Allow jagged rows' are unchecked.

The screenshot shows the Google Cloud Platform console interface. The left sidebar lists various services: Compute Engine, Cloud 기능 (expanded), 네트워킹, Bigtable, Datastore, Storage (selected), SQL, Spanner, STACKDRIVER (expanded), 모니터링, 디버그, 추적, and 로그 기록. The main content area is titled 'Cloud Storage' and '버킷'. It contains a brief description of Cloud Storage and two buttons: '버킷 만들기' (Create Bucket) and '빠른 시작 사용' (Use Quick Start). The browser address bar shows the URL <https://console.cloud.google.com/storage/browser?project=konlper-168808>. The taskbar at the bottom shows several pinned icons.

버킷 만들기

The screenshot shows the Google Cloud Platform Storage bucket creation interface. The URL in the browser is <https://console.cloud.google.com/storage/create-bucket?project=konlper-168808>. The page title is "버킷 만들기" (Create Bucket). The left sidebar shows "Storage" selected. The main form has the following fields:

- 이름:** yarling-stratum-4912
- 기본 저장소 클래스:** Multi-Regional (selected)
- Multi-Regional 위치:** 미국 (selected)
- 라벨 지정:** (empty)

At the bottom are "만들기" (Create) and "취소" (Cancel) buttons.

73/92

The screenshot shows the Google Cloud Platform Storage Bucket creation interface. The left sidebar lists 'Storage' under the 'Storage' category. The main area is titled 'Bucket 만들기' (Create Bucket). The 'Name' field contains 'yarling-stratum-4912'. The 'Basic storage class' dropdown is set to 'Multi-Regional'. Other options shown are 'Regional', 'Nearline', and 'Coldline'. The 'Multi-Regional location' dropdown is set to 'Asia'. At the bottom, there are 'Create' and 'Cancel' buttons.

버킷 확인

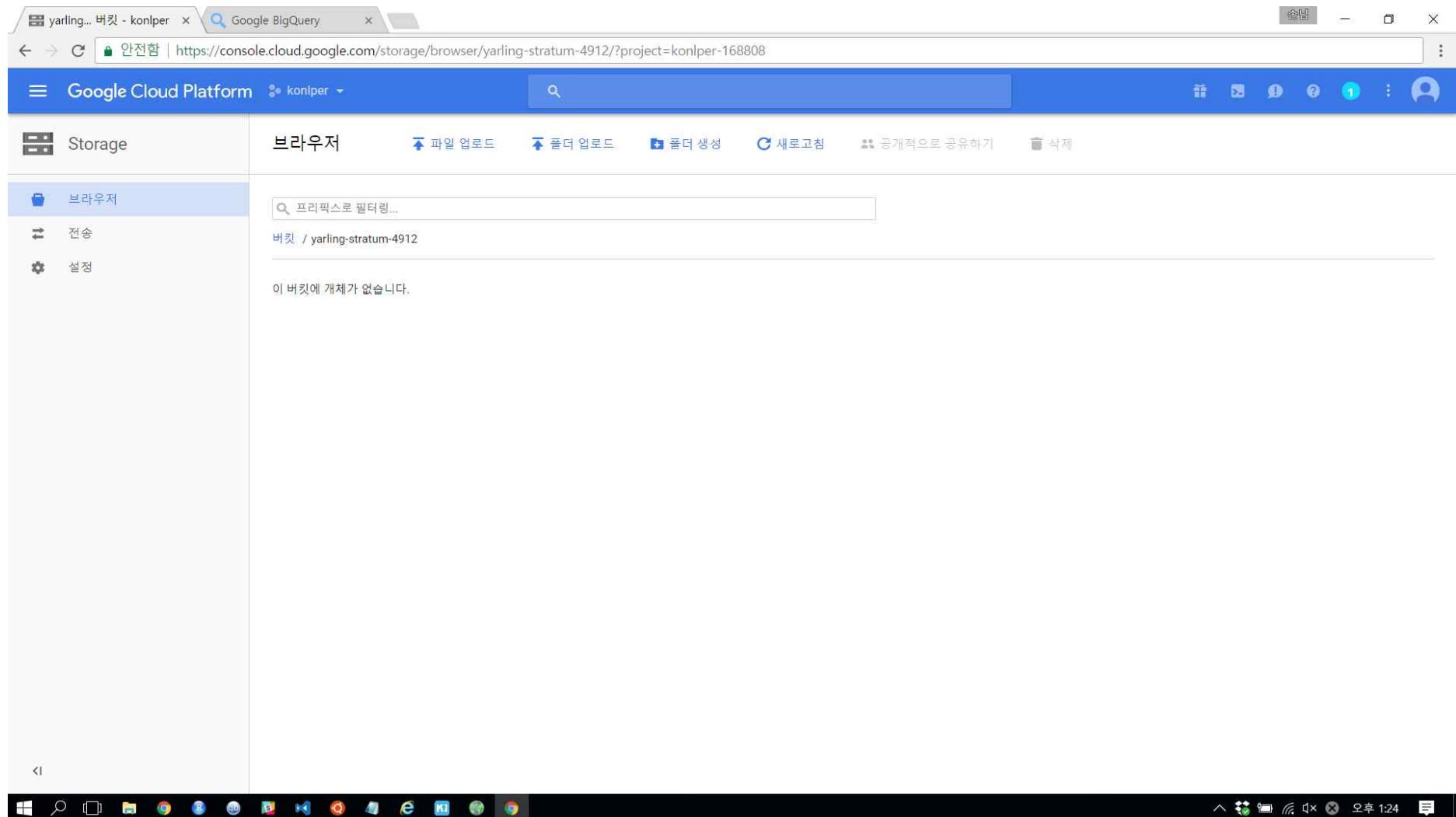
The screenshot shows the Google Cloud Platform Storage browser interface. The left sidebar has 'Storage' selected, with '브라우저' (Browser) highlighted. The main area displays a table with one row of data:

이름	기본 저장소 클래스	위치	라ベル
yarling-stratum-4912	Multi-Regional	ASIA	[More]

At the bottom, the taskbar shows various pinned icons including File Explorer, Task View, and several browser windows.

75/92

폴더 업로드



업로드 확인

The screenshot shows a web browser window for the Google Cloud Platform Storage browser. The URL in the address bar is <https://console.cloud.google.com/storage/browser/yarling-stratum-4912/?project=konlper-168808>. The left sidebar shows 'Storage' and three options: '브라우저' (selected), '전송', and '설정'. The main area displays a list of buckets under '브라우저'. A search bar at the top says '프리픽스로 필터링...'. Below it, a table lists a single folder: 'recomen/'.

이름	크기	유형	저장소 클래스	최종 수정 시간	공개적으로 공유하기
recomen/	-	폴더	-	-	

77/92

업로드 확인

The screenshot shows the Google Cloud Platform Storage browser interface. The left sidebar has 'Storage' selected, with '브라우저' (Browser) highlighted. The main area displays a list of uploaded CSV files in the 'recomen' bucket. The table includes columns for Name, Size, Type, Location, Last Modified, and Share Link.

이름	크기	유형	저장소 클래스	최종 수정 시간	공개적으로 공유하기
chennel.csv	184.7KB	application/vnd.ms-excel	Multi-Regional	17. 6. 28. 오후 1:24	<input type="checkbox"/>
competitor.csv	522.51KB	application/vnd.ms-excel	Multi-Regional	17. 6. 28. 오후 1:24	<input type="checkbox"/>
customer.csv	454.13KB	application/vnd.ms-excel	Multi-Regional	17. 6. 28. 오후 1:24	<input type="checkbox"/>
item.csv	192.9KB	application/vnd.ms-excel	Multi-Regional	17. 6. 28. 오후 1:24	<input type="checkbox"/>
membership.csv	177.42KB	application/vnd.ms-excel	Multi-Regional	17. 6. 28. 오후 1:24	<input type="checkbox"/>

gs:// 경로확인

The screenshot shows the 'Create Table' page in the Google BigQuery web interface. The left sidebar shows 'Query History' and 'Job History' sections, and a 'Public Datasets' section listing various datasets like 'bigquery-public-data:hacker_news'. The main form is titled 'Create Table' and contains the following fields:

- Source Data:** 'Create from source' is selected. A 'Repeat job' dropdown is set to 'Select Previous Job'. The 'Location' is set to 'Google Cloud Storage' with the path 'gs://yarling-stratum-4912/recomen/chennel.csv'. The 'File format' is set to 'CSV'.
- Destination Table:** 'Table name' is set to 'recom . chennel'. 'Table type' is set to 'Native table'.
- Schema:** 'Automatically detect' is checked. A note says 'Schema will be automatically generated.'
- Options:** 'Number of errors allowed' is set to '0'. 'Write preference' is set to 'Write if empty'.

At the bottom left is a blue 'Create Table' button. The top of the window shows browser tabs for 'yarling... 버킷 - konlper', 'Google BigQuery', and 'yarling... 버킷', along with a URL 'https://bigquery.cloud.google.com/repeatloadjob/konlper-168808:bquijob_56cf83cd_15cecf4e97'. The bottom of the screen shows a Windows taskbar with various icons and the time '오후 1:34'.

작업 결과 확인

The screenshot shows the Google BigQuery web interface. The top navigation bar includes tabs for 'yarling... 버킷 - konlper' and 'Google BigQuery'. The address bar shows the URL: <https://bigquery.cloud.google.com/jobs/konlper-168808>. The main content area is titled 'Recent Jobs' and displays two completed load operations:

Job Status	Action	Details	Last Updated
Load (Green)	Repeat Load Job	gs://yarling-stratum-4912/recomen/chennel.csv to konlper-168808:recom.chennel	1:34PM
Load (Red)	Repeat Load Job	gs://yarling-stratum-4912/recomen/chennel.csv to konlper-168808:recom.chennel	1:34PM

The left sidebar features a 'COMPOSE QUERY' button, a 'Query History' section, and a 'Job History' section. Under 'Job History', there is a dropdown menu for filtering by ID or label, currently set to 'konlper'. Below this are sections for 'recom' and 'chennel'. A 'Public Datasets' section lists various datasets from sources like 'bigquery-public-data' and 'gdelt-bq'.

80/92

query 실행화면

The screenshot shows the Google BigQuery web interface. On the left, there's a sidebar with 'COMPOSE QUERY' and sections for 'Query History' and 'Job History'. Below that is a dropdown for 'konlper' with options 'recom' and 'chennel'. Under 'Public Datasets', several datasets are listed: 'bigquery-public-data:hacker_news', 'bigquery-public-data:noaa_gsod', 'bigquery-public-data:samples', 'bigquery-public-data:usa_names', 'gdelt-bq:hathitrustbooks', 'gdelt-bq:internetarchivebooks', 'lookerdata:cdc', 'nyc-tlc:green', and 'nyc-tlc:yellow'. The main area is titled 'New Query' and contains the following SQL code:

```
1 SELECT * FROM [konlper-168808:recom.chennel] LIMIT 100
```

A green status bar below the code says 'Valid: This query will process 259 KB when run.' Below the code are buttons for 'RUN QUERY', 'Save Query', 'Save View', 'Format Query', and 'Show Options'. A checkmark icon is next to the 'RUN QUERY' button. At the bottom, tabs for 'Results', 'Explanation', and 'Job Information' are visible. A red error message 'Query Failed' is displayed, followed by the error details: 'Error: Encountered " "SELECT" "SELECT "" at line 1, column 1. Was expecting: <EOF>' and 'Job ID: konlper-168808:bquijob_3ada8983_15cecfbea1b'. A note at the bottom says 'Note: You can also find errors in your queries before running them. Click the ! below the query composition box to enable real-time validation.' The bottom of the screen shows a Windows taskbar with various icons.

비용 발생 경고

The screenshot shows the Google BigQuery web interface. On the left, there's a sidebar with 'COMPOSE QUERY' and sections for 'Query History' and 'Job History'. Below that is a dropdown menu for 'konlper' with options like 'recom' and 'chennel'. Under 'Public Datasets', several datasets are listed, including 'bigquery-public-data:hacker_news', 'bigquery-public-data:noaa_gsod', 'bigquery-public-data:samples', 'bigquery-public-data:usa_names', 'gdelt-bq:hathitrustbooks', 'gdelt-bq:internetarchivebooks', 'lookerdata:cdc', 'nyc-flc:green', and 'nyc-lfc:yellow'. The main area shows a 'New Query' window with the following SQL code:

```
SELECT * FROM [konlper-168808:recom.chennel] LIMIT 100
```

Below the code, a message says 'Valid: This query will process 259 rows'. There are 'RUN QUERY' and 'Save Query' buttons. A 'Confirm query' dialog box is overlaid on the screen, containing the following text:

Confirm query

With this query, you will be billed for all the data in the table (even if your query contains a LIMIT clause). If you're using the free tier, this query still counts against your free quota.

You can use table preview instead to see records for free and without affecting quotas.

Don't show this again

Buttons: Run query, Go to table preview, Cancel

Note: You can also find errors in your queries before running them. Click the ! below the query composition box to enable real-time validation.



82/92

query 결과 확인

The screenshot shows the Google BigQuery web interface. On the left, there's a sidebar with 'COMPOSE QUERY' buttons for 'Query History' and 'Job History', and a dropdown for 'konlper' which includes 'recom' and 'chennel'. Below that is a list of 'Public Datasets' with links to various datasets like 'bigquery-public-data:hacker_news', 'bigquery-public-data:noaa_gsod', etc. The main area is titled 'New Query' and contains a SQL code editor with the following query:

```
1 SELECT * FROM [konlper-168808:recom.chennel] LIMIT 100
```

A green status bar below the query says 'Valid: This query will process 259 KB when run.' Below the editor are buttons for 'RUN QUERY', 'Save Query', 'Save View', 'Format Query', and 'Show Options'. A message indicates 'Query complete (2.1s elapsed, 259 KB processed)' with a checkmark icon. The results section shows a table with columns 'Row', 'cusID', 'chennel', and 'useCnt'. The data is as follows:

Row	cusID	chennel	useCnt
1	14	A_MOBILE/APP	1
2	74	A_MOBILE/APP	1
3	241	A_MOBILE/APP	1
4	304	A_MOBILE/APP	1
5	326	A_MOBILE/APP	1
6	377	A_MOBILE/APP	1
7	448	A_MOBILE/APP	1
8	518	A_MOBILE/APP	1
9	549	A_MOBILE/APP	1

At the bottom, there are buttons for 'Table' and 'JSON', and a footer navigation bar with links like First, <Prev, Rows 1 - 9 of 100, Next>, Last. The system tray at the bottom right shows the date and time as '오후 1:36'.

공개 데이터셋

The screenshot shows a browser window with multiple tabs open, including 'yarling... 버킷 - konlper', 'Google BigQuery', 'yarling... 버킷', 'Google Cloud Platform', and 'Google BigQuery Public'. The main content is the 'Google BigQuery Public Datasets' documentation page. The page has a sidebar on the left with links to 'All Resources', 'Pricing and Quotas', 'Release Notes', 'Support', and 'Public Datasets' (which is currently selected). The main content area displays five public datasets: '1000 Cannabis Genomes Project', 'Bay Area Bike Share Trips', 'Chicago Crime Data', 'Chicago Taxi Trips', and 'EPA Historical Air Quality Data'. A right sidebar lists more datasets like 'GDELT Book Corpus', 'GitHub Data', and 'Hacker News'. The bottom of the screen shows a Windows taskbar with various icons.

Google BigQuery Public Datasets

A public dataset is any dataset that is stored in BigQuery and made available to the general public. This page lists a special group of public datasets that Google BigQuery hosts for you to access and integrate into your applications. Google pays for the storage of these data sets and provides public access to the data via BigQuery. You pay only for the queries that you perform on the data (the first 1 TB per month is free, subject to [query pricing details](#)).

Public datasets hosted by BigQuery

- 1000 Cannabis Genomes Project**
Genomic open dataset of approximately 850 strains of Cannabis via the Open Cannabis Project.
- Bay Area Bike Share Trips**
This data includes all Bay Area Bike Share trips from August 2013 to the present, and is updated daily.
- Chicago Crime Data**
This dataset reflects reported incidents of crime that occurred in the City of Chicago from 2001 to the present.

목차

- Public datasets hosted by BigQuery
 - 1000 Cannabis Genomes Project
 - Bay Area Bike Share Trips
 - Chicago Crime Data
 - Chicago Taxi Trips
 - EPA Historical Air Quality Data
 - GDELT Book Corpus
 - GitHub Data
 - Hacker News
 - Healthcare Common Procedure Coding System (HCPCS) Level II
 - IRS Form 990 Data
 - Major League Baseball Data
 - Medicare Data
 - NHTSA Traffic Fatality Data
 - NOAA GHCN Weather
 - NOAA GSOD Weather
 - NOAA ICOADS
 - NYC 311 Service Requests
 - NYC Citi Bike Trips
 - NYC TLC Trips

wikipedia dataset

The screenshot shows the Google BigQuery web interface. On the left, there's a sidebar with 'COMPOSE QUERY' and sections for 'Query History' and 'Job History'. Below that is a filter for 'Filter by ID or label' with 'konlper' selected. Under 'konlper', there are two entries: 'recom' and 'chennel'. In the main area, there are tabs for 'Google BigQuery' and 'Google Cloud Platform'. A search bar at the top has 'yarl... 버킷 - konlper' and 'yarl... 버킷' selected. The URL in the address bar is <https://bigquery.cloud.google.com/table/bigquery-public-data:samples.wikipedia>. The main content area shows a 'New Query' window with the following SQL code:

```
SELECT * FROM [konlper-168808:recom.channel] LIMIT 100
```

Below the query is a 'RUN QUERY' button. To the right of the query window, a 'Create Dataset' dialog box is open. It contains fields for 'Dataset ID' (set to 'test'), 'Data location' (set to '(unspecified)'), and 'Data expiration' (set to 'In 5 days'). There are 'OK' and 'Cancel' buttons at the bottom of the dialog. The background shows a table schema for the 'samples.wikipedia' dataset with columns: id (INTEGER, NULLABLE), language (STRING, REQUIRED), wp_namespace (INTEGER, REQUIRED), and is_redirect (BOOLEAN, NULLABLE). The 'language' column is described as being empty in the current dataset. The 'wp_namespace' column is described as being used for namespaces like 'Talk' and 'User'. The 'is_redirect' column is described as being present in versions later than ca. 200908. At the bottom of the interface, there are buttons for 'Query Table', 'Copy Table', 'Export Table', and 'Delete Table'. The status bar at the bottom right shows the time as '오후 1:42'.

dataset 복사

Google BigQuery

New Query

```
1 SELECT * FROM [konlper-168808:recom.channel] LIMIT 100
```

Valid: This query will process 259 KB when run.

RUN QUERY Save Query Save View Format Query Show Options Query complete (2.1s elapsed, 259 KB processed)

Table Details: wikipedia

Schema Details Preview

title	STRING	REQUIRED	The title of the page, as displayed on the page (not in the URL). Always starts with a capital letter and may begin with a namespace (e.g. "Talk:", "User:", "User Talk:", ...)
id	INTEGER	NULLABLE	A unique ID for the article that was revised. These correspond to the order in which articles were created, except for the first several thousand IDs, which are issued in alphabetical order.
language	STRING	REQUIRED	Empty in the current dataset.
is_redirect	BOOLEAN	NULLABLE	Versions later than ca. 200908 may have a redirection marker in the XML.

Copy table Export table Delete table

INT32 REQUIRED Wikipedia segments its pages into namespaces (e.g. "Talk", "User", etc.)
MEDIA = 202; // =2 in WP XML, but these values must be >0
SPECIAL = 201; // =1 in WP XML, but these values must be >0
MAIN = 0;
TALK = 1;
USER = 2;
USER_TALK = 3;

86/92

The screenshot shows the Google BigQuery web interface. On the left, there's a sidebar with 'COMPOSE QUERY' and sections for 'Query History' and 'Job History'. Below that is a search bar labeled 'Filter by ID or label' with 'konlper' selected. Under 'konlper', there are three collapsed categories: 'recom', 'chennel', and 'test'. To the right of the sidebar, a 'New Query' tab is open with the following SQL query:

```
1 SELECT * FROM [konlper-168808:recom.channel] LIMIT 100
```

Below the query editor, there's a 'Copy Table' dialog box. The dialog has fields for 'Destination project' (set to 'konlper (konlper-168808)'), 'Destination dataset' (set to 'test'), and 'Destination table' (set to 'wikipedia_copy'). At the bottom of the dialog are 'OK' and 'Cancel' buttons. The background shows a table schema for the 'wikipedia' dataset, with columns: title (INT64), id (INTEGER, NULLABLE), language (STRING, REQUIRED), wp_namespace (INTEGER, REQUIRED), and is_redirect (BOOLEAN, NULLABLE). The table description indicates that 'id' is a unique ID for articles, 'language' is the language of the article, 'wp_namespace' is the namespace, and 'is_redirect' is a boolean indicating if it's a redirect. The bottom of the screen shows a Windows taskbar with various icons and the system tray.

데이터 구조 보기

Google BigQuery

New Query

```
1 SELECT * FROM [konlper-168808:test.wikipedia_copy] LIMIT 100
```

Valid: This query will process 35.7 GB when run.

RUN QUERY Save Query Save View Format Query Show Options

contributor_id	INTEGER	NULLABLE	Typically, either _id and _username or _ip will be set. A (very) small fraction of edits have neither _ip or (_id and _username). They show up on Wikipedia as "(Username or IP removed)".
contributor_username	STRING	NULLABLE	Typically, either _id and _username or _ip will be set. A (very) small fraction of edits have neither _ip or (_id and _username). They show up on Wikipedia as "(Username or IP removed)".
timestamp	INTEGER	REQUIRED	In Unix time, seconds since epoch.
is_minor	BOOLEAN	NULLABLE	Corresponds to the "Minor Edit" checkbox on Wikipedia's edit page.
is_bot	BOOLEAN	NULLABLE	A special flag that some of Wikipedia's more active bots voluntarily set.
reversion_id	INTEGER	NULLABLE	If this edit is a reversion to a previous edit, this field records the revision_id that was reverted to. If the same article text occurred multiple times, then this will point to the earliest revision. Only revisions with greater than fifty characters are considered for this field. This is to avoid labeling multiple blankings as reversions.
comment	STRING	NULLABLE	Optional user-supplied description of the edit. Section edits are, by default, prefixed with /* Section Name */.
num_characters	INTEGER	REQUIRED	The length of the article after the revision was applied.

Add New Fields

큰 데이터 query 진행

```
select title,sum(num_characters) as num_characters  
from [konlper-168808:test.wikipedia_copy]  
where regexp_match(title,'[Ss]eoul')  
group by title  
order by num_characters desc;
```

query 결과

Google BigQuery

New Query

```
1 select title,sum(num_characters) as num_characters
2 from [konlper-168808:test.wikipedia_copy]
3 where regexp_match(title,'[Ss]eoul')
4 group by title
5
6 order by num_characters desc;
```

Valid: This query will process 9.13 GB when run.

RUN QUERY Save Query Save View Format Query Show Options Query complete (3.3s elapsed, 9.13 GB processed)

Row	title	num_characters
1	Seoul	104656429
2	User talk:JohnnySeoul	78435528
3	FC Seoul	44944412
4	Seoul National University	18422979
5	Talk:Seoul/Archive2	9263684
6	Seoul Subway Line 1	9074261
7	Seoul Foreign School	7576768
8	Seoul Metropolitan Subway	5946441
9	Seoul International School	4911143

Table JSON First <Prev Rows 1 - 9 of 980 Next> Last

90/92

과제

1. RSQLite와 DBI를 활용해서 nycflights13 데이터를 db Table로 만들고, 제출해 주세요.
 1. 'nycflights13' 패키지를 설치하고 5개 데이터를 확인하세요.
 2. dbConnect 명령으로 SQLite 파일을 sql_[이름].db 으로 생성하세요.
 3. 5개의 데이터를 각각의 이름으로 table을 생성하세요.
 4. 1)~3)의 과정을 모두 코드로 남기고 run_[이름].R로 저장하세요.
 5. run_[이름].R과 sql_[이름].db 두 개의 파일을 class2assignment 폴더에 저장하세요.
 6. github에서 pull request로 제출해 주세요.
 7. recomen 폴더에 있는 데이터 6개도 같은 과정을 진행하고, 제출은 하지 마세요.
 8. .gitignore를 이용하면 push의 범위에서 제외할 수 있습니다.

과제

1. recomen 폴더에 있는 데이터 6개를 bigquery에 업로드 하고 질의를 실행해 보세요.
 1. "bigrquery" 패키지의 기능으로 작은 용량 5개의 데이터를 업로드해 주세요.
 2. Storage 서비스를 이용해서 tran.csv 파일을 bigquery에 테이블로 생성하세요.
 3. query_exec 함수와 "select * from [tran] limit 10" 을 실행하고 결과를 받으세요.
 - [tran]은 각자 해당하는 이름으로 변경하셔야 합니다.
 4. 1)~3)의 과정을 big_[이름].R로 저장하세요.
 5. big_[이름].R 파일을 class2assignment 폴더에 저장하세요.
 6. github에서 pull request로 제출해 주세요.