

COMS W4111-002, V02 (Spring 2022)

Introduction to Databases

Homework 2: Programming and Non-Programming

Due Wednesday, February 23, 2022 at 11:59 PM

Introduction

Overview

This notebook has **2 sections that you must complete**:

1. Written questions testing knowledge of concepts. Answering these questions may require reviewing lecture slides, slides associated with the textbook, and/or online material. Both tracks complete this section.
1. Practical problems involving data modeling, relational algebra and SQL. Both tracks complete this section.

We will separately release the track-specific Programming and Non-Programming parts of HW2.

Submission

You will **submit 2 files** for this assignment.

1. Submit a zip file titled `<your_uni>_hw2_all.zip` to **HW2 All - Zip** on Gradescope.
 - Replace `<your_uni>` with your uni. My submission would be titled `dff9_hw2_all.zip`.
 - The zipped directory you submit should contain the following files:
 - `<your_uni>_hw2_all.ipynb`
 - `Appearances.csv`
 - `Batting.csv`
 - `People.csv`
 - Any image files you choose to embed in your notebook.
 - All of these files, except the images you may embed in your notebook, are included in `s22_w4111_hw2_all.zip`, which you downloaded from Courseworks. You will have to rename the notebook file you downloaded to `<your_uni>_hw2_all.ipynb`, as discussed above.

1. Submit a PDF file titled `<your_uni>_hw2_all.pdf` to **HW2 All - PDF** on Gradescope.

- This should be a PDF of your completed HW2 All Python notebook.
- **Tag pages for each problem.** Per course policy, any untagged submission will receive an automatic 0.
- Double check your submission on Gradescope to ensure that the PDF conversion worked and that your pages are appropriately tagged.

Collaboration and Information

- Answering some of the questions may require independent research to find information. We encourage you to try troubleshooting problems independently before reaching out for help.
- You may use any information you get in TA or Prof. Ferguson's office hours, from lectures or from recitations. This includes slides related to the recommended textbook.
- You may use information that you find on the web.
- You are NOT allowed to collaborate with other students outside of office hours.

Written Questions

Question 1: NULL

Briefly explain Codd's 3rd Rule.

- What are some interpretations of a NULL value?
- An alternative to using NULL is some other value for indicating missing data, e.g. using -1 for the value of a weight column. Explain the benefits of NULL relative to other approaches.

Answer:

Codd's 3rd rule is the systematic treatment of null values. This means that NULL values in a database must be given a systematic and uniform treatment.

A NULL value can be interpreted any of the following – data that is missing, data that is not known, or data that is not applicable.

The benefits of using NULL relative to other approaches is that the NULL value prevents programmers from having to carefully code some SQL statements to match the specific, varying choices they made. Without the NULL value, users must understand the domain to determine "unknown" values or know what choice a developer made. Essentially, the NULL value is used as a standard and does not require any additional programming or memorization since NULL is pretty well-known and accepted. It allows the DBMS to be fully functional and behave in a predictable and standard manner, irrespective of what data type of a value is missing/unknown/inapplicable.

Question 2: Keys

Briefly explain the following concepts:

- Primary Key
- Candidate Key
- Super Key
- Alternate Key
- Composite Key
- Unique Key
- Foreign Key

Answer:

Primary Key: used to denote a candidate key that is chosen by the database designer as the principal means of identifying entities within an entity set. Primary key is an attribute of the table that is unique and doesn't contain any NULL values; a primary key helps to identify entities of a table.

Candidate Key: are a subset of super keys for which no proper subset is a superkey (minimal super key). Like a super key, a candidate key also identifies each tuple in a table uniquely. A candidate key's attribute can accept NULL value.

Super Key: a set of one or more attributes that, taken collectively, allow us to identify uniquely an entity in the entity set. A super key can contain redundant attributes that might not be important for identifying tuples.

Alternate Key: are those candidate keys that were not chosen to be the table's primary key.

Composite Key: is a candidate key [DOUBLE CHECK THIS] that consists of more than one attribute. A composite key is a combination of two or more columns: sometimes no single attribute will have the property to uniquely identify tuples in a table, thus we use a group of attributes to ensure uniqueness.

Unique Key: is a column or set of columns that uniquely identify each database record. All values have to be unique in this key and a unique key can have only one null value.

Foreign Key: column or group of columns used to provide a link between data in two tables. A foreign key may accept non-unique and NULL values. Foreign keys acts as a cross-reference between tables by linking up to a primary key in a separate table.

Question 3: Algebra

Briefly explain what it means for the relational algebra to be *closed* under the operations in the algebra. What is an important benefit?

Answer:

What it means for relational algebra to be closed under the operations in the algebra is that both the operands and outputs are relations. The query language takes an instance of relations

and does operations that work on one or more relations to describe another relations without altering the original relation(s). So output from one operation can turn into the input of another operation, which allows expressions to be nested in the relational algebra – this is an important benefit.

Source: <https://www.w3schools.in/dbms/relational-algebra/>

Question 4: Equivalent Queries

Briefly explain the concept of equivalent queries. Use the concept to explain how it is possible to derive the JOIN operation from other operations (SELECT, PROJECT).

Answer:

In SQL, equivalent queries are queries that produce the same output given any input relation.

Using this concept, we can derive JOIN operation from other operations: The JOIN operator pairs two tuples from different relations, if and only if a given condition is satisfied. Thus, the JOIN operation is a combination of a Cartesian product followed by a selection process.

Consider the following example: there are two tables, instructor and teaches, and we want to get only those tuples that pertain to instructors and the courses that they taught. We can accomplish this by doing instructor join teaches. However, another equivalent way is to get the cartesian product of instructor x teaches and then select those tuples which pertain to instructors and the courses they taught, using selection. This is done by selecting tuples where instructor.id = teaches.id on instructor x teaches. The result of these two methods are equivalent.

Question 5: More General Attribute Types

The relational model places restrictions on attributes. Many data scenarios have more complex types of attributes. **Briefly** explain the following types of attributes:

- Simple attribute
- Composite attribute
- Derived attribute
- Single-value attribute
- Multi-value attribute

Answer:

Simple attribute: is not divided into subparts, and is, thus, referred to as atomic values since you cannot divide it further.

Example: "firstName".

Composite attribute: can be divided into subparts (i.e. other attributes). Composite attributes are good to use if a user plans to refer to an entire attribute on some occasions, and to only a

component of the attribute on other occasions.

Example: "fullName" which can be divided into "firstName", "middleName", and "lastName".

Derived attribute: is an attribute whose value can be derived from the values of other related attributes or entities. The value of a derived attribute is not stored, but is computed when required.

Example: Someone's age can be derived by the attribute "date-of-birth" and calculating its difference from the current date.

Single-value attribute: have a single value for a particular entry.

Example: Date of birth – a person can have only one.

Multi-value attribute: have a set of values for a specific entry. Sometimes, when appropriate, upper and lower bounds may be placed on the number of values in a multi-valued attribute.

Example: Consider an employee entity set with the attribute "phone-number". An employee may have zero, one, or several phone numbers, and different employees may have different numbers of phones. Thus, "phone-number" is multi-valued.

Practical Problems

Setup

- Modify the cells below to setup your environment.
- The change should just be setting the DB user ID and password, replacing my user ID and password with yours for MySQL.

```
In [54]: database_user_id = "root"
         database_pwd = "dbuserdbuser"
```

```
In [55]: database_url = "mysql+pymysql://" + \
         database_user_id + ":" + database_pwd + "@localhost"
         database_url
```

```
Out[55]: 'mysql+pymysql://root:dbuserdbuser@localhost'
```

```
In [56]: %reload_ext sql
```

```
In [57]: %sql $database_url
```

```
Out[57]: 'Connected: root@None'
```

```
In [58]: from sqlalchemy import create_engine
```

```
In [59]: sqla_engine = create_engine(database_url)
```

```
In [60]: #
# We are going to create a schema and some tables for the HW.
#
%sql create schema if not exists S22_W4111_HW2
%sql select 1;

* mysql+pymysql://root:***@localhost
1 rows affected.
* mysql+pymysql://root:***@localhost
1 rows affected.
```

```
Out[60]: 1
1
```

Question 6: Manipulating String and Types

Setup

- Run the following code cells.
- These cells create a table `people_info` and loads the table with a bunch of input strings.

```
In [61]: input_string = [
    "Towny,Cavet,tcavet0@blinklist.com,1/9/1971,+62 (340) 387-5141",
    "Port,Gaylor,pgaylor1@blogger.com,3/15/1939,+86 (517) 758-9970",
    "Georgetta,Haddon,ghaddon2@symantec.com,9/19/1997,+81 (356) 753-5556",
    "Wylma,Lanney,wlanney3@list-manage.com,2/21/2018,+385 (853) 541-7347",
    "Mignonne,Georgeson,mgeorgeson4@123-reg.co.uk,8/7/1991,+63 (834) 397-5285",
    "Cchaddie,Cossins,ccossins5@chronoengine.com,3/12/1911,+242 (313) 943-4080",
    "Andie,Matyushonok,amatyushonok6@ask.com,4/24/1907,+380 (410) 464-9093",
    "Skiptie,Zuenelli,szuenelli7@merriam-webster.com,3/22/2014,+7 (279) 484-2088",
    "Averyl,Barajas,abarajas8@fastcompany.com,6/19/1996,+232 (962) 344-7325",
    "Olia,Habens,ohabens9@quantcast.com,2/28/1922,+98 (935) 300-9359"
]
```

```
In [62]: import pandas
```

```
In [63]: df = pandas.DataFrame(input_string)
```

```
In [64]: df.to_sql(
    "people_info", con=sqla_engine, if_exists="replace", index=False,
    schema="S22_W4111_HW2")
```

```
In [65]: %sql use S22_W4111_HW2
```

```
%sql select 1;
```

```
* mysql+pymysql://root:***@localhost
0 rows affected.
* mysql+pymysql://root:***@localhost
1 rows affected.
```

```
Out[65]: 1
1
```

- Test loading the data.

```
In [66]: %sql select * from people_info
```

```
* mysql+pymysql://root:***@localhost
10 rows affected.
```

```
Out[66]: 0
```

```
Towny,Cavet,tcavet0@blinklist.com,1/9/1971,+62 (340) 387-5141
Port,Gaylor,pgaylor1@blogger.com,3/15/1939,+86 (517) 758-9970
Georgetta,Haddon,ghaddon2@symantec.com,9/19/1997,+81 (356) 753-5556
Wylma,Lanney,wlanney3@list-manage.com,2/21/2018,+385 (853) 541-7347
Mignonne,Georgeson,mgeorgeson4@123-reg.co.uk,8/7/1991,+63 (834) 397-5285
Cchaddie,Cossins,ccossins5@chronoengine.com,3/12/1911,+242 (313) 943-4080
Andie,Matyushonok,amatyushonok6@ask.com,4/24/1907,+380 (410) 464-9093
Skippie,Zuenelli,szuenelli7@merriam-webster.com,3/22/2014,+7 (279) 484-2088
Averyl,Barajas,abarajas8@fastcompany.com,6/19/1996,+232 (962) 344-7325
Olia,Habens,ohabens9@quantcast.com,2/28/1922,+98 (935) 300-9359
```

- Can we describe what the table looks like?

```
In [67]: %sql describe people_info;
```

```
* mysql+pymysql://root:***@localhost
1 rows affected.
```

```
Out[67]: Field Type Null Key Default Extra
0 text YES None
```

Tasks

- The created table has one column `0` . The values are strings with data separated by `,` . The fields in the string are (in order):
 - `first_name`
 - `last_name`
 - `email`

- `date_of_birth`
 - `telephone_no`, which is of the form `+CC (XXX)–XXX–XXXX` where `CC` is the country code and the remainder is the number.
- You must process and cleanup the data using **ONLY** SQL statements. The cleanup tasks include:
 - Creating a new table `people_info_clean` with a structure that better represents the data, e.g. columns, column data types, etc.
 - Converting each string and its subfields into the rows of `people_info_clean`.
 - You may use as many DDL and DML SQL statements as you need.
 - Execute your statements in the cells below and show the output of the execution.
 - The last two cells show show the data and schema for the information.

In [68]:

```
%%sql

drop table if exists people_info_clean;

create table if not exists S22_W4111_HW2.people_info_clean (
    first_name VARCHAR(100) null,
    last_name VARCHAR(100) null,
    email VARCHAR(100) not null primary key,
    date_of_birth date null,
    telephone_no VARCHAR(100) null
);

select * from people_info_clean;
```

```
* mysql+pymysql://root:***@localhost
0 rows affected.
0 rows affected.
0 rows affected.
```

Out[68]: first_name last_name email date_of_birth telephone_no

In [69]:

```
%%sql
insert into people_info_clean (first_name, last_name, email,
                              date_of_birth, telephone_no)
select
    substring_index(`0`, ",", 1)
      as first_name, #first_name
    substring_index(substring_index(`0`, ",", 2), ",", -1)
      as last_name, #last name
    substring_index(substring_index(`0`, ",", 3), ",", -1)
      as email, #email
    str_to_date(substring_index(substring_index(`0`, ",", 4), ",", -1),
               "%c/%e/%Y") as date_of_birth, #date_of_birth
    substring_index(`0`, ",", -1)
      as telephone_no #telephone_no
from
    people_info
;
```



```
* mysql+pymysql://root:***@localhost
10 rows affected.
```

Out[69]: []

In [70]:

```
#
# Show the data.
#
%sql select * from people_info_clean
```

```
* mysql+pymysql://root:***@localhost
10 rows affected.
```

Out[70]:

first_name	last_name	email	date_of_birth	telephone_no
Averyl	Barajas	abarajas8@fastcompany.com	1996-06-19	+232 (962) 344-7325
Andie	Matyushonok	amatyushonok6@ask.com	1907-04-24	+380 (410) 464-9093
Cchaddie	Cossins	ccossins5@chronoengine.com	1911-03-12	+242 (313) 943-4080
Georgetta	Haddon	ghaddon2@symantec.com	1997-09-19	+81 (356) 753-5556
Mignonne	Georgeson	mgeorgeson4@123-reg.co.uk	1991-08-07	+63 (834) 397-5285
Olia	Habens	ohabens9@quantcast.com	1922-02-28	+98 (935) 300-9359
Port	Gaylor	pgaylor1@blogger.com	1939-03-15	+86 (517) 758-9970
Skippie	Zuenelli	szuenelli7@merriam-webster.com	2014-03-22	+7 (279) 484-2088
Towny	Cavet	tcavet0@blinklist.com	1971-01-09	+62 (340) 387-5141
Wylma	Lanney	wlanney3@list-manage.com	2018-02-21	+385 (853) 541-7347

In [71]:

```
#
# Show the schema (architecture and structure).
#
%sql describe people_info_clean;
```

```
* mysql+pymysql://root:***@localhost
5 rows affected.
```

Out[71]:

	Field	Type	Null	Key	Default	Extra
	first_name	varchar(100)	YES		None	
	last_name	varchar(100)	YES		None	
	email	varchar(100)	NO	PRI	None	
	date_of_birth	date	YES		None	
	telephone_no	varchar(100)	YES		None	

Question 7: Intermediate SQL and Data Processing

Task 1: Load Data

- Continue to use the schema you created - S22_W4111_HW2 .
- There are three files in the homework folder:

- People.csv
 - Appearances.csv
 - Batting.csv
- Use one of the approaches we have previously used directly in notebooks to load the CSV files into the schema above.
 - You **may not** use external tools like DataGrip.
 - Some examples of techniques are in HW 1 and in the Pandas examples.
 - Put your code in the cells provided below. The final cells, which you must run after loading the CSV files, simply display some information.

In [72]:

Your code

In [73]:

```

%%sql
use S22_W4111_HW2;

drop table if exists people;

create table if not exists S22_W4111_HW2.people (
    playerID VARCHAR(100) null,
    birthYear int null,
    birthMonth int null,
    birthDay int null,
    birthCountry VARCHAR(100) null,
    birthState VARCHAR(100) null,
    birthCity VARCHAR(100) null,
    deathYear int null,
    deathMonth int null,
    deathDay int null,
    deathCountry VARCHAR(100) null,
    deathState VARCHAR(100) null,
    deathCity VARCHAR(100) null,
    nameFirst VARCHAR(100) null,
    nameLast VARCHAR(100) null,
    nameGiven VARCHAR(100) null,
    weight int null,
    height int null,
    bats VARCHAR(100) null,
    throws VARCHAR(100) null,
    debut VARCHAR(100) null,
    finalGame VARCHAR(100) null,
    retroID VARCHAR(100) null,
    bbrefID VARCHAR(100) null
);

select * from people;

```

```

* mysql+pymysql://root:***@localhost
0 rows affected.
0 rows affected.
0 rows affected.
0 rows affected.

```

Out[73]: playerID birthYear birthMonth birthDay birthCountry birthState birthCity deathYear deathMc

In [74]:

```

%%sql drop table if exists S22_W4111_HW2.appearances;

create table if not exists S22_W4111_HW2.appearances (
    yearID int null,
    teamID VARCHAR(100) null,
    lgID VARCHAR(100) null,
    playerID VARCHAR(100) null,
    G_all int null,
    GS int null,
    G_batting int null,
    G_defense int null,
    G_p int null,
    G_c int null,
    G_1b int null,
    G_2b int null,
    G_3b int null,
    G_ss int null,
    G_if int null,
    G_cf int null,
    G_rf int null,
    G_of int null,
    G_dh int null,
    G_ph int null,
    G_pr int null
);

select * from S22_W4111_HW2.appearances;

```

```

* mysql+pymysql://root:***@localhost
0 rows affected.
0 rows affected.
0 rows affected.

```

Out[74]: **yearID teamID lgID playerID G_all GS G_batting G_defense G_p G_c G_1b G_2b G_3b**

In [75]:

```

%%sql drop table if exists S22_W4111_HW2.batting;

create table if not exists S22_W4111_HW2.batting (
    playerID VARCHAR(100) null,
    yearID int null,
    stint int null,
    teamID VARCHAR(100) null,
    lgID VARCHAR(100) null,
    G int null,
    AB int null,
    R int null,
    H int null,
    2B int null,
    3B int null,
    HR int null,
    RBI int null,
    SB int null,
    CS int null,
    BB int null,
    SO int null,
    IBB int null,
    HBP int null,
    SH int null,

```

```

        SF int null,
        GDP int null
    );

select * from S22_W4111_HW2.batting;

```

```

* mysql+pymysql://root:***@localhost
0 rows affected.
0 rows affected.
0 rows affected.

```

Out[75]:

playerID	yearID	stint	teamID	IgID	G	AB	R	H	2B	3B	HR	RBI	SB	CS	BB	SO	IBB	HB
----------	--------	-------	--------	------	---	----	---	---	----	----	----	-----	----	----	----	----	-----	----

```

In [76]: project_root = "/Users/anyadevgan/Desktop/Databases/s22_w4111_hw2_all"

people_df = pandas.read_csv(project_root + "/People.csv")
people_df.to_sql("people", con=sqldb_engine, if_exists="replace",
                  index=False, schema="S22_W4111_HW2")

appearance_df = pandas.read_csv(project_root + "/Appearances.csv")
appearance_df.to_sql("appearances", con=sqldb_engine, if_exists="replace",
                    index=False, schema="S22_W4111_HW2")

batting_df = pandas.read_csv(project_root + "/Batting.csv")
batting_df.to_sql("batting", con=sqldb_engine, if_exists="replace",
                  index=False, schema="S22_W4111_HW2")

```

```

In [77]: # Some tests

```

```

In [78]: %sql select * from people limit 10;

```

```

* mysql+pymysql://root:***@localhost
10 rows affected.

```

Out[78]:

playerID	birthYear	birthMonth	birthDay	birthCountry	birthState	birthCity	deathYear	deathMonth	deathDay
aardsda01	1981.0	12.0	27.0	USA	CO	Denver	None		
aaronha01	1934.0	2.0	5.0	USA	AL	Mobile	2021.0		
aaronto01	1939.0	8.0	5.0	USA	AL	Mobile	1984.0		
aasedo01	1954.0	9.0	8.0	USA	CA	Orange	None		
abadan01	1972.0	8.0	25.0	USA	FL	Palm Beach	None		
abadfe01	1985.0	12.0	17.0	D.R.	La Romana	La Romana	None		
abadijo01	1850.0	11.0	4.0	USA	PA	Philadelphia	1905.0		

playerID	birthYear	birthMonth	birthDay	birthCountry	birthState	birthCity	deathYear	deathMonth	deathDay
abbated01	1877.0	4.0	15.0	USA	PA	Latrobe	1957.0		
abbeybe01	1869.0	11.0	11.0	USA	VT	Essex	1962.0		
abbeych01	1866.0	10.0	14.0	USA	NE	Falls City	1926.0		

In [79]:

```
%sql select * from appearances limit 10;
```

```
* mysql+pymysql://root:***@localhost
10 rows affected.
```

Out[79]:

yearID	teamID	lgID	playerID	G_all	GS	G_batting	G_defense	G_p	G_c	G_1b	G_2b	G_3b
1871	TRO	None	abercda01	1	1.0	1	1.0	0	0	0	0	
1871	RC1	None	addybo01	25	25.0	25	25.0	0	0	0	22	
1871	CL1	None	allisar01	29	29.0	29	29.0	0	0	0	2	
1871	WS3	None	allisdo01	27	27.0	27	27.0	0	27	0	0	
1871	RC1	None	ansonca01	25	25.0	25	25.0	0	5	1	2	
1871	FW1	None	armstbo01	12	12.0	12	12.0	0	0	0	0	
1871	RC1	None	barkeal01	1	1.0	1	1.0	0	0	0	0	
1871	BS1	None	barnero01	31	31.0	31	31.0	0	0	0	16	
1871	FW1	None	barrebi01	1	1.0	1	1.0	0	1	0	0	
1871	BS1	None	barrofr01	18	17.0	18	18.0	0	0	0	1	

In [80]:

```
%sql select * from batting limit 10;
```

```
* mysql+pymysql://root:***@localhost
10 rows affected.
```

Out[80]:

playerID	yearID	stint	teamID	lgID	G	AB	R	H	2B	3B	HR	RBI	SB	CS	BB	SO
abercda01	1871	1	TRO	None	1	4	0	0	0	0	0	0.0	0.0	0.0	0	0.0
addybo01	1871	1	RC1	None	25	118	30	32	6	0	0	13.0	8.0	1.0	4	0.0
allisar01	1871	1	CL1	None	29	137	28	40	4	5	0	19.0	3.0	1.0	2	5.0
allisdo01	1871	1	WS3	None	27	133	28	44	10	2	2	27.0	1.0	1.0	0	2.0
ansonca01	1871	1	RC1	None	25	120	29	39	11	3	0	16.0	6.0	2.0	2	1.0
armstbo01	1871	1	FW1	None	12	49	9	11	2	1	0	5.0	0.0	1.0	0	1.0
barkeal01	1871	1	RC1	None	1	4	0	1	0	0	0	2.0	0.0	0.0	1	0.0
barnero01	1871	1	BS1	None	31	157	66	63	10	9	0	34.0	11.0	6.0	13	1.0
barrebi01	1871	1	FW1	None	1	5	1	1	1	0	0	1.0	0.0	0.0	0	0.0
barrofr01	1871	1	BS1	None	18	86	13	13	2	1	0	11.0	1.0	0.0	0	0.0

In [81]: `%sql describe people;`

```
* mysql+pymysql://root:***@localhost
24 rows affected.
```

Out[81]:

Field	Type	Null	Key	Default	Extra
-------	------	------	-----	---------	-------

playerID	text	YES		None	
birthYear	double	YES		None	
birthMonth	double	YES		None	
birthDay	double	YES		None	
birthCountry	text	YES		None	
birthState	text	YES		None	
birthCity	text	YES		None	
deathYear	double	YES		None	
deathMonth	double	YES		None	
deathDay	double	YES		None	
deathCountry	text	YES		None	
deathState	text	YES		None	
deathCity	text	YES		None	
nameFirst	text	YES		None	
nameLast	text	YES		None	
nameGiven	text	YES		None	
weight	double	YES		None	
height	double	YES		None	
bats	text	YES		None	
throws	text	YES		None	
debut	text	YES		None	
finalGame	text	YES		None	
retroID	text	YES		None	
bbrefID	text	YES		None	

In [82]: `%sql describe appearances;`

```
* mysql+pymysql://root:***@localhost
21 rows affected.
```

Out[82]:

Field	Type	Null	Key	Default	Extra
-------	------	------	-----	---------	-------

yearID	bigint	YES		None	
teamID	text	YES		None	
lgID	text	YES		None	
playerID	text	YES		None	

Field	Type	Null	Key	Default	Extra
G_all	bigint	YES		None	
GS	double	YES		None	
G_batting	bigint	YES		None	
G_defense	double	YES		None	
G_p	bigint	YES		None	
G_c	bigint	YES		None	
G_1b	bigint	YES		None	
G_2b	bigint	YES		None	
G_3b	bigint	YES		None	
G_ss	bigint	YES		None	
G_lf	bigint	YES		None	
G_cf	bigint	YES		None	
G_rf	bigint	YES		None	
G_of	bigint	YES		None	
G_dh	double	YES		None	
G_ph	double	YES		None	
G_pr	double	YES		None	

In [83]:

```
%sql describe batting;
```

```
* mysql+pymysql://root:***@localhost
22 rows affected.
```

Out[83]:

Field	Type	Null	Key	Default	Extra
playerID	text	YES		None	
yearID	bigint	YES		None	
stint	bigint	YES		None	
teamID	text	YES		None	
lgID	text	YES		None	
G	bigint	YES		None	
AB	bigint	YES		None	
R	bigint	YES		None	
H	bigint	YES		None	
2B	bigint	YES		None	
3B	bigint	YES		None	
HR	bigint	YES		None	
RBI	double	YES		None	

Field	Type	Null	Key	Default	Extra
SB	double	YES		None	
CS	double	YES		None	
BB	bigint	YES		None	
SO	double	YES		None	
IBB	double	YES		None	
HBP	double	YES		None	
SH	double	YES		None	
SF	double	YES		None	
GIDP	double	YES		None	

Task 2: Complicated Queries

Note: Performing the query in this task may require changing column values or types.

Query - Career Summary

- Write a query that produces a result of the form:
 - playerID
 - nameLast
 - nameFirst
 - The sum of `appearances.G_all` for the player over all rows.
 - The sum over all rows of the following columns from batting:
 - G
 - AB
 - R
 - H
 - 2B
 - 3B
 - HR
 - RBI
 - BB
 - batting_average , which is defined as $\frac{sum(H)}{sum(AB)}$
 - on_base_percentage , which is defined as $\frac{(sum(H) + sum(BB))}{(sum(AB) + sum(BB))}$
- The query should be limited to 20 rows, and sorted by `on_base_percentage` from highest to lowest.
- `batting_average` and `on_base_percentage` should round to three decimal places.


```
In [84]: %%sql
with people_test as
(
    select playerID, nameLast, nameFirst from people
),
appearances_test as
(
    select playerID, SUM(G_all) as G_all_sum from appearances group by playerID
),
people_appearances as
(
    select * from people_test join appearances_test using (playerID)
),
batting_test as
(
    select playerID, SUM(G) as G_sum,
           SUM(AB) as AB_sum,
           SUM(R) as R_sum,
           SUM(2B) as 2B_sum,
           SUM(3B) as 3B_sum,
           SUM(HR) as HR_sum,
           SUM(RBI) as RBI_sum,
           SUM(BB) as BB_sum,
           ROUND(SUM(H)/if(SUM(AB)=0,NULL,SUM(AB)), 3) as batting_average,
           ROUND((SUM(H)+SUM(BB))/if((sum(AB)+sum(BB))=0,NULL,
           (sum(AB)+sum(BB))), 3) as on_base_percentage
    from batting group by playerID
),
people_appearances_batting as
(
    select * from people_appearances join batting_test using (playerID)
)
select * from people_appearances_batting order by on_base_percentage desc limit
```

```
* mysql+pymysql://root:***@localhost
```

```
20 rows affected.
```

```
Out[84]:
```

playerID	nameLast	nameFirst	G_all_sum	G_sum	AB_sum	R_sum	2B_sum	3B_sum	HR_
torrejo02	Torres	Jose	66	66	1	1	0	0	
meansjo01	Means	John	42	42	1	0	0	0	
sotogr01	Soto	Gregory	60	60	2	0	0	0	
alanirj01	Alaniz	R. J.	12	12	1	0	0	0	
carsoro01	Carson	Robert	31	31	0	1	0	0	
horstje01	Horst	Jeremy	72	72	1	0	0	0	
thompaa01	Thompson	Aaron	52	52	0	0	0	0	
alberan01	Albers	Andrew	26	26	1	0	0	0	
meekev01	Meek	Evan	179	179	1	0	0	0	
melanma01	Melancon	Mark	606	606	0	0	0	0	
hoovejj01	Hoover	J. J.	290	290	0	1	0	0	
schlibr01	Schlitter	Brian	84	84	1	0	0	0	
tupmama01	Tupman	Matt	1	1	1	0	0	0	

playerID	nameLast	nameFirst	G_all_sum	G_sum	AB_sum	R_sum	2B_sum	3B_sum	HR_
montgst01	Montgomery	Steve	72	72	1	1	0	0	
hancory01	Hancock	Ryan	11	11	1	1	0	0	
cammaer01	Cammack	Eric	8	8	1	0	0	1	
parrajo01	Parra	Jose	82	82	0	0	0	0	
duvalmi01	Duvall	Mike	53	53	0	0	0	0	
yanes01	Yan	Esteban	472	472	2	1	0	0	
bruneju01	Brunette	Justin	4	4	1	0	0	0	

Question 8: "Fun" with Sets

- `People` represents basic information about people associated with Major League Baseball.
- `Appearances` contains information about people who appeared (played in) MLB games.
- There are some entries in the `People` table that do not appear in `Appearances`.
- Using a **subquery**, write a query that counts the number of people in the `People` table who do not have an entry in `Appearances`.
- Run your query below. Note, your query will be **SLOW**.

In [88]:

```
%%sql
use S22_W4111_HW2;

SELECT count(playerID)
from people
where playerID not in (select playerID from appearances where playerID = people.
limit 20;
```

```
* mysql+pymysql://root:***@localhost
0 rows affected.
1 rows affected.
```

Out [88]:

```
count(playerID)
460
```

- Just for the heck of it, run the scripts below and repeat your query. Also, the changes I am making are a good hint on how to solve the problem.

In [89]:

```
%%sql
use s22_w4111_hw2;

drop table if exists people_fast;
drop table if exists appearances_fast;
```

```

create table people_fast as select * from people;
create table appearances_fast as select * from appearances;

ALTER TABLE `appearances_fast`
CHANGE COLUMN `playerID` `playerID` VARCHAR(16) NULL DEFAULT NULL ,
ADD INDEX `playerID_idx` (playerID) VISIBLE;

ALTER TABLE `people_fast`
CHANGE COLUMN `playerID` `playerID` VARCHAR(16) NULL DEFAULT NULL ,
ADD INDEX `peopleID_idx` (playerID) VISIBLE;

```

```

* mysql+pymysql://root:***@localhost
0 rows affected.
0 rows affected.
0 rows affected.
20358 rows affected.
108717 rows affected.
108717 rows affected.
20358 rows affected.

```

Out[89]: []

- Run your query here.

In [83]:

```

%%sql use S22_W4111_HW2;

SELECT count(playerID)
from people
where playerID not in (select playerID from appearances
                       where playerID = people.playerID)
limit 20;

```

```

* mysql+pymysql://root:***@localhost
0 rows affected.
1 rows affected.

```

Out[83]: count(playerID)

460

In []:

Question 9: Don Plays Baseball

- I always wanted to play baseball for the Boston Red Sox, and also play with Ted Williams.
- Ted Williams' playerID is willlite01 .
- My playerID would be fergusdo .
- Perform the following tasks using SQL:
 - Insert an entry in people with:
 - playerID = fergusdo

- nameLast = Ferguson
- nameFirst = Donald
- Existence without Ted Williams is meaningless. So, using an Update statement, update the entry in people for fergusdo to have the same birthYear, birthMonth and birthDay as Ted Williams.
- Run a query showing the row in people for fergusdo.
- Delete the row you added.

```
In [90]: # Insert statement
#
%sql insert into people(playerID, nameLast, nameFirst) values("fergusdo", "Fergu

* mysql+pymysql://root:***@localhost
1 rows affected.
```

Out[90]: []

```
In [91]: %%sql
# Update statement
#

update people p1, (select * from people where playerID = "willite01") as p2
set p1.birthYear = p2.birthYear,
p1.birthMonth = p2.birthMonth,
p1.birthDay = p2.birthDay
where p1.playerID = "fergusdo" ;

* mysql+pymysql://root:***@localhost
1 rows affected.
```

Out[91]: []

```
In [92]: # Select statement showing row. not showing updated values
#
%sql select * from people where playerID = "fergusdo";

* mysql+pymysql://root:***@localhost
1 rows affected.
```

Out[92]:

playerID	birthYear	birthMonth	birthDay	birthCountry	birthState	birthCity	deathYear	deathMc
fergusdo	1918.0	8.0	30.0	None	None	None	None	N

```
In [93]: # Delete the created row.
#
%sql delete from people where playerID = "fergusdo";

* mysql+pymysql://root:***@localhost
1 rows affected.
```

Out[93]: []

In [94]:

```
%sql select * from people where playerID = "fergusdo";
```

```
* mysql+pymysql://root:***@localhost
```

```
0 rows affected.
```

```
Out[94]: playerID birthYear birthMonth birthDay birthCountry birthState birthCity deathYear deathMc
```

Question 10: There is No Question 10

- You all get a free point for putting up with me.

```
In [ ]:
```