

You don't have to follow the above workflow to get full credit on assignments, but we suggest it to reduce debugging headaches. We've created the scratchwork `%%sql` cells for you in this assignment, but **do not** add cells between this `%%sql` cell and the Python cell right below it. It will cause errors when we run the autograder, and it will sometimes cause a failure to generate the PDF file.

### 3 Part 1: The IMDb (mini) Dataset

Let's explore a miniature version of the [IMDb Dataset](#). This is the same dataset that we will use for the upcoming homework.

Let's load in the database in two ways (using both Python and cell magic) so that we can flexibly explore the database.

```
[10]: engine = sqlalchemy.create_engine("sqlite:///imdbmini.db")
      connection = engine.connect()
```

```
[11]: %%sql sqlite:///imdbmini.db
```

```
[12]: %%sql
      SELECT * FROM sqlite_master WHERE type='table';
```

```
* sqlite:///imdbmini.db
Done.
```

```
[12]: [('table', 'Title', 'Title', 2, 'CREATE TABLE "Title" (\n"tconst" INTEGER,\n"titleType" TEXT,\n "primaryTitle" TEXT,\n "originalTitle" TEXT,\n "isAdult" TEXT,\n "startYear" TEXT,\n "endYear" TEXT,\n "runtimeMinutes" TEXT,\n "genres" TEXT\n)'),
      ('table', 'Name', 'Name', 12, 'CREATE TABLE "Name" (\n"nconst" INTEGER,\n"primaryName" TEXT,\n "birthYear" TEXT,\n "deathYear" TEXT,\n "primaryProfession" TEXT\n)'),
      ('table', 'Role', 'Role', 70, 'CREATE TABLE "Role" (\n"tconst" INTEGER,\n"ordering" INTEGER,\n"category" TEXT,\n"job" TEXT,\n"characters" TEXT\n)'),
      ('table', 'Rating', 'Rating', 41, 'CREATE TABLE "Rating" (\n"tconst" INTEGER,\n"averageRating" TEXT,\n"numVotes" TEXT\n)')]
```

From running the above cell, we see the database has 4 tables: `Name`, `Role`, `Rating`, and `Title`.

[Click to Expand] See descriptions of each table's schema.

**Name** – Contains the following information for names of people.

- `nconst` (text) - alphanumeric unique identifier of the name/person
- `primaryName` (text)– name by which the person is most often credited
- `birthYear` (integer) – in YYYY format
- `deathYear` (integer) – in YYYY format

**Role** – Contains the principal cast/crew for titles.

- `tconst` (text) - alphanumeric unique identifier of the title
- `ordering` (integer) – a number to uniquely identify rows for a given `tconst`

- nconst (text) - alphanumeric unique identifier of the name/person
- category (text) - the category of job that person was in
- characters (text) - the name of the character played if applicable, else '\N'

**Rating** – Contains the IMDb rating and votes information for titles.

- tconst (text) - alphanumeric unique identifier of the title
- averageRating (real) – weighted average of all the individual user ratings
- numVotes (integer) - number of votes (i.e., ratings) the title has received

**Title** - Contains the following information for titles.

- tconst (text) - alphanumeric unique identifier of the title
- titleType (text) - the type/format of the title
- primaryTitle (text) - the more popular title / the title used by the filmmakers on promotional materials at the point of release
- isAdult (text) - 0: non-adult title; 1: adult title
- year (YYYY) – represents the release year of a title.
- runtimeMinutes (integer) – primary runtime of the title, in minutes

From the above descriptions, we can conclude the following: \* `Name.nconst` and `Title.tconst` are primary keys of the `Name` and `Title` tables, respectively. \* that `Role.nconst` and `Role.tconst` are **foreign keys** that point to `Name.nconst` and `Title.tconst`, respectively.

### 3.1 Question 1

What are the different kinds of `titleTypes` included in the `Title` table? Write a query to find out all the unique `titleTypes` of films using the `DISTINCT` keyword. (**You may not use `GROUP BY`.**)

```
[19]: %%sql
/*
 * Code in this scratchwork cell is __not graded.__
 * Copy over any SQL queries you write here into the below Python cell.
 * Do __not__ insert any new cells in between the SQL/Python cells!
 * Doing so may break the autograder.
 */
-- Write below this comment. --
SELECT DISTINCT titleType
FROM Title;
```

```
* sqlite:///imdbmini.db
Done.
```

```
[19]: [('short',),
      ('movie',),
      ('tvSeries',),
      ('tvMovie',),
      ('tvMiniSeries',),
      ('video',),
      ('videoGame',),
```

```
('tvEpisode',),  
( 'tvSpecial',)]
```

```
[21]: query_q1 = """  
SELECT DISTINCT titleType  
FROM Title;  
"""  
  
res_q1 = pd.read_sql(query_q1, engine)  
res_q1
```

```
[21]:      titleType  
0      short  
1      movie  
2    tvSeries  
3    tvMovie  
4 tvMiniSeries  
5      video  
6  videoGame  
7    tvEpisode  
8    tvSpecial
```

```
[22]: grader.check("q1")
```

```
[22]: q1 results: All test cases passed!
```

---

## 3.2 Question 2

Before we proceed we want to get a better picture of the kinds of jobs that exist. To do this examine the `Role` table by computing the number of records with each job `category`. Present the results in descending order by the total counts.

The top of your table should look like this (however, you should have more rows):

	category	total
<b>0</b>	actor	21665
<b>1</b>	writer	13830
<b>2</b>	...	...

```
[25]: %%sql  
/*  
 * Code in this scratchwork cell is __not graded.__  
 * Copy over any SQL queries you write here into the below Python cell.  
 * Do __not__ insert any new cells in between the SQL/Python cells!  
 * Doing so may break the autograder.
```

```

*/
-- Write below this comment. --
SELECT category, COUNT(*) AS total
FROM Role
GROUP BY category
ORDER BY total DESC;

```

```

* sqlite:///imdbmini.db
Done.

```

```

[25]: [('actor', 21665),
      ('writer', 13830),
      ('actress', 12175),
      ('producer', 11028),
      ('director', 6995),
      ('composer', 4123),
      ('cinematographer', 2747),
      ('editor', 1558),
      ('self', 623),
      ('production_designer', 410),
      ('archive_footage', 66),
      ('archive_sound', 6)]

```

```

[27]: query_q2 = """
SELECT category, COUNT(*) AS total
FROM Role
GROUP BY category
ORDER BY total DESC;
"""

res_q2 = pd.read_sql(query_q2, engine)
res_q2

```

```

[27]:
      category  total
0         actor  21665
1         writer  13830
2        actress  12175
3        producer  11028
4        director   6995
5         composer   4123
6  cinematographer  2747
7          editor   1558
8           self    623
9  production_designer   410
10    archive_footage    66
11    archive_sound     6

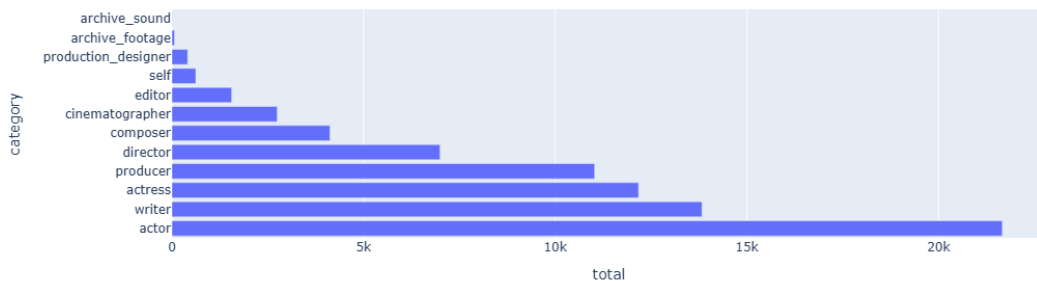
```

```
[28]: grader.check("q2")
```

[28]: q2 results: All test cases passed!

If we computed the results correctly we should see a nice horizontal bar chart of the counts per category below:

```
[29]: # just run this cell
px.bar(res_q2, x="total", y="category", orientation='h')
```



### 3.3 Question 3

Now that we have a better sense of the basics of our data, we can ask some more interesting questions.

The `Rating` table has the `numVotes` and the `averageRating` for each title. Which 10 films have the most ratings?

Write a SQL query that outputs three fields: the `title`, `numVotes`, and `averageRating` for the 10 films that have the highest number of ratings. Sort the result in descending order by the number of votes.

*Hint:* The `numVotes` in the `Rating` table is not an integer! Use `CAST(Rating.numVotes AS int) AS numVotes` to convert the attribute to an integer.

```
[32]: """sql
/*
 * Code in this scratchwork cell is __not graded.__
 * Copy over any SQL queries you write here into the below Python cell.
 * Do __not__ insert any new cells in between the SQL/Python cells!
 * Doing so may break the autograder.
 */
-- Write below this comment. --
SELECT primaryTitle AS title, CAST(Rating.numVotes AS int) AS numVotes,
       averageRating
FROM Rating
```

```
JOIN Title
ON Rating.tconst = Title.tconst
ORDER BY numVotes DESC
LIMIT 10;
```

```
* sqlite:///imdbmini.db
Done.
```

```
[32]: [('The Shawshank Redemption', 2462686, '9.3'),
      ('The Dark Knight', 2417875, '9.0'),
      ('Inception', 2169255, '8.8'),
      ('Fight Club', 1939312, '8.8'),
      ('Pulp Fiction', 1907561, '8.9'),
      ('Forrest Gump', 1903969, '8.8'),
      ('Game of Thrones', 1874040, '9.2'),
      ('The Matrix', 1756469, '8.7'),
      ('The Lord of the Rings: The Fellowship of the Ring', 1730296, '8.8'),
      ('The Lord of the Rings: The Return of the King', 1709023, '8.9')]
```

```
[33]: query_q3 = """
      SELECT primaryTitle AS title, CAST(Rating.numVotes AS int) AS numVotes,
             averageRating
      FROM Rating
      JOIN Title
      ON Rating.tconst = Title.tconst
      ORDER BY numVotes DESC
      LIMIT 10;
      """

      res_q3 = pd.read_sql(query_q3, engine)
      res_q3
```

```
[33]:
```

	title	numVotes	averageRating
0	The Shawshank Redemption	2462686	9.3
1	The Dark Knight	2417875	9.0
2	Inception	2169255	8.8
3	Fight Club	1939312	8.8
4	Pulp Fiction	1907561	8.9
5	Forrest Gump	1903969	8.8
6	Game of Thrones	1874040	9.2
7	The Matrix	1756469	8.7
8	The Lord of the Rings: The Fellowship of the Ring	1730296	8.8
9	The Lord of the Rings: The Return of the King	1709023	8.9

```
[34]: grader.check("q3")
```

[34]: q3 results: All test cases passed!

## 4 Part 2: Election Donations in New York City

Finally, let's analyze the Federal Election Commission (FEC)'s public records. We connect to the database in two ways (using both Python and cell magic) so that we can flexibly explore the database.

```
[35]: # important!!! run this cell and the next one
import sqlalchemy
# create a SQL Alchemy connection to the database
engine = sqlalchemy.create_engine("sqlite:///fec_nyc.db")
connection = engine.connect()
```

```
[36]: %sql sqlite:///fec_nyc.db
```

### 4.1 Table Descriptions

Run the below cell to explore the **schemas** of all tables saved in the database.

If you'd like, you can consult the below linked FEC pages for the descriptions of the tables themselves.

- `cand` ([link](#)): Candidates table. Contains names and party affiliation.
- `comm` ([link](#)): Committees table. Contains committee names and types.
- `indiv_sample_nyc` ([link](#)): All individual contributions from New York City .

```
[37]: %%sql
/* just run this cell */
SELECT sql FROM sqlite_master WHERE type='table';
```

```
* sqlite:///fec_nyc.db
  sqlite:///imdbmini.db
Done.
```

```
[37]: [('CREATE TABLE "cand" (\n    cand_id character varying(9),\n    cand_name\n    text,\n    cand_pty_affiliation character varying(3),\n    cand_election_yr i\n    ... (196 characters truncated) ... er varying(9),\n    cand_st1 text,\n    cand_st2 text,\n    cand_city text,\n    cand_st character varying(2),\n    cand_zip character varying(10)\n)\n'),\n ('CREATE TABLE "comm"(\n    "cmte_id" TEXT,\n    "cmte_nm" TEXT,\n    "tres_nm"\n    TEXT,\n    "cmte_st1" TEXT,\n    "cmte_st2" TEXT,\n    "cmte_city" TEXT,\n    "cmte_s\n    ... (46 characters truncated) ... XT,\n    "cmte_tp" TEXT,\n    "cmte_pty_affiliation" TEXT,\n    "cmte_filing_freq" TEXT,\n    "org_tp" TEXT,\n    "connected_org_nm" TEXT,\n    "cand_id" TEXT\n)\n'),\n ('CREATE TABLE indiv_sample_nyc (\n    cmte_id character varying(9),\n    amndt_ind character(1),\n    rpt_tp character varying(3),\n    transaction_pg\n    ... (299 characters truncated) ... transaction_amt integer,\n    other_id
```

```
text,\n    tran_id text,\n    file_num bigint,\n    memo_cd text,\n    memo_text text,\n    sub_id bigint\n)',,)]
```

Let's look at the `indiv_sample_nyc` table. The below cell displays individual donations made by residents of the state of New York. We use `LIMIT 5` to avoid loading and displaying a huge table.

```
[38]: %%sql
/* just run this cell */
SELECT *
FROM indiv_sample_nyc
LIMIT 5;
```

```
* sqlite:///fec_nyc.db
sqlite:///imdbmini.db
Done.
```

```
[38]: [('C00445015', 'N', 'Q1', 'P', 15951128130, '15', 'IND', 'SINGER, TRIPP MR.', 'NEW YORK', 'NY', '100214505', 'ATLANTIC MAILBOXES, INC.', 'OWNER', '01302015', 1000, '', 'A-CF13736', 1002485, '', '', 4041420151241812398),
('C00510461', 'N', 'Q1', 'P', 15951129284, '15E', 'IND', 'SIMON, DANIEL A', 'NEW YORK', 'NY', '100237940', 'N/A', 'RETIRED', '03292015', 400, 'C00401224', 'VN8JBDDJBA8', 1002590, '', '* EARMARKED CONTRIBUTION: SEE BELOW', 4041420151241813640),
('C00422410', 'N', 'Q1', 'P', 15970352211, '15', 'IND', 'ABDUL RAUF, FEISAL', 'NEW YORK', 'NY', '101150010', 'THE CORDOBA INITIATIVE', 'CHAIRMAN', '03042015', 250, '', 'VN8A3DBSYG6', 1003643, '', '', 4041620151241914560),
('C00510461', 'N', 'Q1', 'P', 15951129280, '15', 'IND', 'SCHWARZER, FRANK', 'NEW YORK', 'NY', '100145135', 'METRO HYDRAULIC JACK CO', 'SALES', '01162015', 100, '', 'VN8JBDAP4C4', 1002590, '', '* EARMARKED CONTRIBUTION: SEE BELOW', 4041420151241813630),
('C00510461', 'N', 'Q1', 'P', 15951129281, '15', 'IND', 'SCHWARZER, FRANK', 'NEW YORK', 'NY', '100145135', 'METRO HYDRAULIC JACK CO', 'SALES', '02162015', 100, '', 'VN8JBDBRDG3', 1002590, '', '* EARMARKED CONTRIBUTION: SEE BELOW', 4041420151241813632)]
```

You can write a SQL query to return the id and name of the first five candidates from the Democratic party, as below:

```
[39]: %%sql
/* just run this cell */
SELECT cand_id, cand_name
FROM cand
WHERE cand_pty_affiliation = 'DEM'
LIMIT 5;
```

```
* sqlite:///fec_nyc.db
sqlite:///imdbmini.db
Done.
```



```
[39]: [('HOAL05049', 'CRAMER, ROBERT E "BUD" JR'),
      ('HOAL07086', 'SEWELL, TERRY CINA ANDREA'),
      ('HOAL07094', 'HILLIARD, EARL FREDERICK JR'),
      ('HOAR01091', 'GREGORY, JAMES CHRISTOPHER'),
      ('HOAR01109', 'CAUSEY, CHAD')]
```

## 4.2 [Tutorial] Matching Text with LIKE

First, let's look at 2016 election contributions made by Donald Trump, who was a New York (NY) resident during that year. The following SQL query returns the `cmte_id`, `transaction_amt`, and `name` for every contribution made by any donor with “DONALD” and “TRUMP” in their name in the `indiv_sample_nyc` table.

Notes: \* We use the `WHERE ... LIKE '...'` to match fields with text patterns. The `%` wildcard represents at least zero characters. Compare this to what you know from regex! \* We use `pd.read_sql` syntax here because we will do some EDA on the result `res`.

```
[40]: # just run this cell
example_query = """
SELECT
    cmte_id,
    transaction_amt,
    name
FROM indiv_sample_nyc
WHERE name LIKE '%TRUMP%' AND name LIKE '%DONALD%';
"""

example_res = pd.read_sql(example_query, engine)
example_res
```

```
[40]:
```

	cmte_id	transaction_amt	name
0	C00230482	2600	DONALD, TRUMP
1	C00230482	2600	DONALD, TRUMP
2	C00014498	9000	TRUMP, DONALD
3	C00494229	2000	TRUMP, DONALD MR
4	C00571869	2700	TRUMP, DONALD J.
..	...	...	...
152	C00608489	5	DONALD J TRUMP FOR PRESIDENT INC
153	C00608489	5	DONALD J TRUMP FOR PRESIDENT INC
154	C00608489	5	DONALD J TRUMP FOR PRESIDENT INC
155	C00608489	5	DONALD J TRUMP FOR PRESIDENT INC
156	C00608489	5	DONALD J TRUMP FOR PRESIDENT INC

```
[157 rows x 3 columns]
```

If we look at the list above, it appears that some donations were not by Donald Trump himself, but instead by an entity called “DONALD J TRUMP FOR PRESIDENT INC”. Fortunately, we see that our query only seems to have picked up one such anomalous name.

```
[41]: # just run this cell
example_res['name'].value_counts()
```

```
[41]: TRUMP, DONALD J.                133
DONALD J TRUMP FOR PRESIDENT INC    15
TRUMP, DONALD                      4
DONALD, TRUMP                      2
TRUMP, DONALD MR                   1
TRUMP, DONALD J MR.                1
TRUMP, DONALD J MR                 1
Name: name, dtype: int64
```

### 4.3 Question 4

Revise the above query so that the 15 anomalous donations made by “DONALD J TRUMP FOR PRESIDENT INC” do not appear. Your resulting table should have 142 rows.

Hints: \* Consider using the above query as a starting point, or checking out the SQL query skeleton at the top of this lab. \* The NOT keyword may also be useful here.

```
[44]: %sql
/*
 * Code in this scratchwork cell is __not graded.__
 * Copy over any SQL queries you write here into the below Python cell.
 * Do __not__ insert any new cells in between the SQL/Python cells!
 * Doing so may break the autograder.
 */
-- Write below this comment. --
SELECT
    cmte_id,
    transaction_amt,
    name
FROM indiv_sample_nyc
WHERE name LIKE '%TRUMP%' AND name LIKE '%DONALD%' AND name NOT LIKE '%INC%';
```

```
* sqlite:///fec_nyc.db
sqlite:///imdbmini.db
Done.
```

```
[44]: [('C00230482', 2600, 'DONALD, TRUMP'),
('C00230482', 2600, 'DONALD, TRUMP'),
('C00014498', 9000, 'TRUMP, DONALD'),
('C00494229', 2000, 'TRUMP, DONALD MR'),
('C00571869', 2700, 'TRUMP, DONALD J.'),
('C00571869', 2700, 'TRUMP, DONALD J.'),
('C00136457', 5000, 'TRUMP, DONALD'),
('C00034033', 5000, 'TRUMP, DONALD'),
('C00580100', 4049, 'TRUMP, DONALD J.'),
```

('C00554949', 2600, 'TRUMP, DONALD J MR.'),  
 ('C00369033', 1000, 'TRUMP, DONALD J MR'),  
 ('C00580100', 514, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 4049, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 4049, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 4958, 'TRUMP, DONALD J.'),  
 ('C00580100', 1713, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 4127, 'TRUMP, DONALD J.'),  
 ('C00580100', 3063, 'TRUMP, DONALD J.'),  
 ('C00580100', 4164, 'TRUMP, DONALD J.'),  
 ('C00580100', 3954, 'TRUMP, DONALD J.'),  
 ('C00580100', 2460, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 2460, 'TRUMP, DONALD J.'),  
 ('C00580100', 4591, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 224, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 2460, 'TRUMP, DONALD J.'),  
 ('C00580100', 2426, 'TRUMP, DONALD J.'),  
 ('C00580100', 3548, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 2529, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 2544, 'TRUMP, DONALD J.'),  
 ('C00580100', 2531, 'TRUMP, DONALD J.'),  
 ('C00580100', 2426, 'TRUMP, DONALD J.'),  
 ('C00580100', 3499, 'TRUMP, DONALD J.'),  
 ('C00580100', 1216, 'TRUMP, DONALD J.'),  
 ('C00580100', 2460, 'TRUMP, DONALD J.'),  
 ('C00580100', 7990, 'TRUMP, DONALD J.'),  
 ('C00580100', 5574, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 2544, 'TRUMP, DONALD J.'),  
 ('C00580100', 6250, 'TRUMP, DONALD J.'),  
 ('C00580100', 2991, 'TRUMP, DONALD J.'),  
 ('C00580100', 2544, 'TRUMP, DONALD J.'),  
 ('C00580100', 4931, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 2544, 'TRUMP, DONALD J.'),

('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 2693, 'TRUMP, DONALD J.'),  
 ('C00580100', 2620, 'TRUMP, DONALD J.'),  
 ('C00580100', 2693, 'TRUMP, DONALD J.'),  
 ('C00580100', 2692, 'TRUMP, DONALD J.'),  
 ('C00580100', 3371, 'TRUMP, DONALD J.'),  
 ('C00580100', 2679, 'TRUMP, DONALD J.'),  
 ('C00580100', 2587, 'TRUMP, DONALD J.'),  
 ('C00580100', 3360, 'TRUMP, DONALD J.'),  
 ('C00580100', 2587, 'TRUMP, DONALD J.'),  
 ('C00580100', 2587, 'TRUMP, DONALD J.'),  
 ('C00580100', 3347, 'TRUMP, DONALD J.'),  
 ('C00580100', 2587, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 2536, 'TRUMP, DONALD J.'),  
 ('C00580100', 2536, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 2536, 'TRUMP, DONALD J.'),  
 ('C00580100', 3706, 'TRUMP, DONALD J.'),  
 ('C00580100', 2536, 'TRUMP, DONALD J.'),  
 ('C00580100', 2536, 'TRUMP, DONALD J.'),  
 ('C00580100', 3706, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 5297, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 3706, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 11023, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),  
 ('C00580100', 3706, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 3706, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00055582', 10000, 'TRUMP, DONALD'),  
 ('C00580100', 9651, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 2000000, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 9106, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 2574, 'TRUMP, DONALD J.'),  
 ('C00580100', 15000, 'TRUMP, DONALD J.'),

```
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 9000, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 6197, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 8735, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 2000000, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 10000000, 'TRUMP, DONALD J.' ),
( 'C00580100', 9000, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 6990, 'TRUMP, DONALD J.' ),
( 'C00580100', 1029, 'TRUMP, DONALD J.' ),
( 'C00580100', 8896, 'TRUMP, DONALD J.' ),
( 'C00580100', 9000, 'TRUMP, DONALD J.' ),
( 'C00580100', 7519, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 9013, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 10418, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 7401, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 2000000, 'TRUMP, DONALD J.' ),
( 'C00580100', 9000, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 9000, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 9752, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' ),
( 'C00580100', 23775, 'TRUMP, DONALD J.' ),
( 'C00580100', 2000000, 'TRUMP, DONALD J.' ),
( 'C00580100', 2574, 'TRUMP, DONALD J.' )]
```

```
[45]: query_q4 = """
SELECT
    cmte_id,
    transaction_amt,
    name
FROM indiv_sample_nyc
WHERE name LIKE '%TRUMP%' AND name LIKE '%DONALD%' AND name NOT LIKE '%INC%';
```

```
"""
```

```
res_q4 = pd.read_sql(query_q4, engine)
res_q4
```

```
[45]:
```

	cmte_id	transaction_amt	name
0	C00230482	2600	DONALD, TRUMP
1	C00230482	2600	DONALD, TRUMP
2	C00014498	9000	TRUMP, DONALD
3	C00494229	2000	TRUMP, DONALD MR
4	C00571869	2700	TRUMP, DONALD J.
..	...	...	...
137	C00580100	9752	TRUMP, DONALD J.
138	C00580100	2574	TRUMP, DONALD J.
139	C00580100	23775	TRUMP, DONALD J.
140	C00580100	2000000	TRUMP, DONALD J.
141	C00580100	2574	TRUMP, DONALD J.

```
[142 rows x 3 columns]
```

```
[46]: grader.check("q4")
```

```
[46]: q4 results: All test cases passed!
```

## 4.4 Question 5: JOINing Tables

Let's explore the other two tables in our database: `cand` and `comm`.

The `cand` table contains summary financial information about each candidate registered with the FEC or appearing on an official state ballot for House, Senate or President.

```
[47]: %%sql
/* just run this cell */
SELECT *
FROM cand
LIMIT 5;
```

```
* sqlite:///fec_nyc.db
* sqlite:///imdbmini.db
```

```
Done.
```

```
[47]: [('HOAK00097', 'COX, JOHN R.', 'REP', 2014, 'AK', 'H', 0, 'C', 'N', 'C00525261',
'P.O. BOX 1092 ', '', 'ANCHOR POINT', 'AK', '99556'),
('HOAL02087', 'ROBY, MARTHA', 'REP', 2016, 'AL', 'H', 2, 'I', 'C', 'C00462143',
'PO BOX 195', '', 'MONTGOMERY', 'AL', '36101'),
('HOAL02095', 'JOHN, ROBERT E JR', 'IND', 2016, 'AL', 'H', 2, 'C', 'N', '',
'1465 W OVERBROOK RD', '', 'MILLBROOK', 'AL', '36054'),
('HOAL05049', 'CRAMER, ROBERT E "BUD" JR', 'DEM', 2008, 'AL', 'H', 5, '', 'P',
```

```
'C00239038', 'PO BOX 2621', '', 'HUNTSVILLE', 'AL', '35804'),
('HOAL05163', 'BROOKS, MO', 'REP', 2016, 'AL', 'H', 5, 'I', 'C', 'C00464149',
'7610 FOXFIRE DRIVE', '', 'HUNTSVILLE', 'AL', '35802')]
```

The `comm` table contains summary financial information about each committee registered with the FEC. Committees are organizations that spend money for political action or parties, or spend money for or against political candidates.

```
[48]: %%sql
/* just run this cell */
SELECT *
FROM comm
LIMIT 5;
```

```
* sqlite:///fec_nyc.db
  sqlite:///imdbmini.db
Done.
```

```
[48]: [('C00000059', 'HALLMARK CARDS PAC', 'ERIN BROWER', '2501 MCGEE', 'MD#288',
'KANSAS CITY', 'MO', '64108', 'U', 'Q', 'UNK', 'M', 'C', '', ''),
('C00000422', 'AMERICAN MEDICAL ASSOCIATION POLITICAL ACTION COMMITTEE',
'WALKER, KEVIN', '25 MASSACHUSETTS AVE, NW', 'SUITE 600', 'WASHINGTON', 'DC',
'20001', 'B', 'Q', '', 'M', 'M', 'AMERICAN MEDICAL ASSOCIATION', ''),
('C00000489', 'D R I V E POLITICAL FUND CHAPTER 886', 'TOM RITTER', '3528 W
RENO', '', 'OKLAHOMA CITY', 'OK', '73107', 'U', 'N', '', 'Q', 'L', 'TEAMSTERS
LOCAL UNION 886', ''),
('C00000547', 'KANSAS MEDICAL SOCIETY POLITICAL ACTION COMMITTEE', 'C. RICHARD
BONEBRAKE, M.D.', '623 SW 10TH AVE', '', 'TOPEKA', 'KS', '66612', 'U', 'Q',
'UNK', 'Q', 'T', '', ''),
('C00000638', 'INDIANA STATE MEDICAL ASSOCIATION POLITICAL ACTION COMMITTEE',
'VIDYA KORA, M.D.', '322 CANAL WALK, CANAL LEVEL', '', 'INDIANAPOLIS', 'IN',
'46202', 'U', 'Q', '', 'Q', 'M', '', '')]
```

#### 4.4.1 Question 5a

Notice that both the `cand` and `comm` tables have a `cand_id` column. Let's try joining these two tables on this column to print out committee information for candidates.

List the first 5 candidate names (`cand_name`) in reverse lexicographic order by `cand_name`, along with their corresponding committee names. **Only select rows that have a matching `cand_id` in both tables.**

Your output should look similar to the following:

	cand_name	cmte_nm
0	ZUTLER, DANIEL PAUL MR	CITIZENS TO ELECT DANIEL P ZUTLER FOR PRESIDENT
1	ZUMWALT, JAMES	ZUMWALT FOR CONGRESS
...	...	...

Consider starting from the following query skeleton, which uses the `AS` keyword to rename the `cand` and `comm` tables to `c1` and `c2`, respectively. Which join is most appropriate?

```
SELECT ...
FROM cand AS c1
    [INNER | {LEFT | RIGHT | FULL } {OUTER}] JOIN comm AS c2
    ON ...
...
...;
```

```
[50]: %%sql
/*
 * Code in this scratchwork cell is __not graded.__
 * Copy over any SQL queries you write here into the below Python cell.
 * Do __not__ insert any new cells in between the SQL/Python cells!
 * Doing so may break the autograder.
 */
-- Write below this comment. --
SELECT cand_name, cmte_nm
FROM cand AS c1
JOIN comm AS c2
ON c1.cand_id = c2.cand_id
ORDER BY cand_name DESC
LIMIT 5;

* sqlite:///fec_nyc.db
* sqlite:///imdbmini.db
Done.
```

```
[50]: [('ZUTLER, DANIEL PAUL MR', 'CITIZENS TO ELECT DANIEL P ZUTLER FOR PRESIDENT'),
      ('ZUMWALT, JAMES', 'ZUMWALT FOR CONGRESS'),
      ('ZUKOWSKI, ANDREW GEORGE', 'ZUKOWSKI FOR CONGRESS'),
      ('ZUCCOLO, JOE', 'JOE ZUCCOLO FOR CONGRESS'),
      ('ZORN, ROBERT ERWIN', 'CONSTITUTIONAL COMMITTEE')]
```

```
[51]: query_q5a = """
SELECT cand_name, cmte_nm
FROM cand AS c1
JOIN comm AS c2
ON c1.cand_id = c2.cand_id
```



```
ORDER BY cand_name DESC
LIMIT 5;
"""
```

```
res_q5a = pd.read_sql(query_q5a, engine)
res_q5a
```

```
[51]:
```

	cand_name	cmte_nm
0	ZUTLER, DANIEL PAUL MR	CITIZENS TO ELECT DANIEL P ZUTLER FOR PRESIDENT
1	ZUMWALT, JAMES	ZUMWALT FOR CONGRESS
2	ZUKOWSKI, ANDREW GEORGE	ZUKOWSKI FOR CONGRESS
3	ZUCCOLO, JOE	JOE ZUCCOLO FOR CONGRESS
4	ZORN, ROBERT ERWIN	CONSTITUTIONAL COMMITTEE

```
[52]: grader.check("q5a")
```

```
[52]: q5a results: All test cases passed!
```

#### 4.4.2 Question 5b

Suppose we modify the query from the previous part to include *all* candidates, **including those that don't have a committee**.

List the first 5 candidate names (`cand_name`) in reverse lexicographic order by `cand_name`, along with their corresponding committee names. If the candidate has no committee in the `comm` table, then `cmte_nm` should be `NULL` (or `None` in the Python representation).

Your output should look similar to the following:

	cand_name	cmte_nm
0	ZUTLER, DANIEL PAUL MR	CITIZENS TO ELECT DANIEL P ZUTLER FOR PRESIDENT
...	...	...
4	ZORNOW, TODD MR	None

Hint: Start from the same query skeleton as the previous part. Which join is most appropriate?

```
[53]: %%sql
/*
 * Code in this scratchwork cell is __not graded.__
 * Copy over any SQL queries you write here into the below Python cell.
 * Do __not__ insert any new cells in between the SQL/Python cells!
 * Doing so may break the autograder.
 */
-- Write below this comment. --
```

```
SELECT cand_name, cmte_nm
FROM cand AS c1
LEFT JOIN comm AS c2
ON c1.cand_id = c2.cand_id
ORDER BY cand_name DESC
LIMIT 5;
```

```
* sqlite:///fec_nyc.db
  sqlite:///imdbmini.db
Done.
```

```
[53]: [('ZUTLER, DANIEL PAUL MR', 'CITIZENS TO ELECT DANIEL P ZUTLER FOR PRESIDENT'),
      ('ZUMWALT, JAMES', 'ZUMWALT FOR CONGRESS'),
      ('ZUKOWSKI, ANDREW GEORGE', 'ZUKOWSKI FOR CONGRESS'),
      ('ZUCCOLO, JOE', 'JOE ZUCCOLO FOR CONGRESS'),
      ('ZORNOW, TODD MR', None)]
```

```
[54]: query_q5b = """
      SELECT cand_name, cmte_nm
      FROM cand AS c1
      LEFT JOIN comm AS c2
      ON c1.cand_id = c2.cand_id
      ORDER BY cand_name DESC
      LIMIT 5;
      """

      res_q5b = pd.read_sql(query_q5b, engine)
      res_q5b
```

```
[54]:
```

	cand_name	cmte_nm
0	ZUTLER, DANIEL PAUL MR	CITIZENS TO ELECT DANIEL P ZUTLER FOR PRESIDENT
1	ZUMWALT, JAMES	ZUMWALT FOR CONGRESS
2	ZUKOWSKI, ANDREW GEORGE	ZUKOWSKI FOR CONGRESS
3	ZUCCOLO, JOE	JOE ZUCCOLO FOR CONGRESS
4	ZORNOW, TODD MR	None

```
[55]: grader.check("q5b")
```

```
[55]: q5b results: All test cases passed!
```

## 4.5 Question 6: Subqueries and Grouping (OPTIONAL)

If we return to our results from Question 4, we see that many of the contributions were to the same committee:

```
[ ]: # Your SQL query result from Question 4
      # reprinted for your convenience
      res_q4['cmte_id'].value_counts()
```

---

Create a new SQL query that returns the total amount that Donald Trump contributed to each committee.

Your table should have four columns: `cmte_id`, `total_amount` (total amount contributed to that committee), `num_donations` (total number of donations), and `cmte_nm` (name of the committee). Your table should be sorted in **decreasing order** of `total_amount`.

**This is a hard question!** Don't be afraid to reference the lecture slides, or the overall SQL query skeleton at the top of this lab.

Here are some other hints:

- Note that committee names are not available in `indiv_sample_nyc`, so you will have to obtain information somehow from the `comm` table (perhaps a JOIN would be useful).
- Remember that you can compute summary statistics after grouping by using aggregates like `COUNT(*)`, `SUM()` as output fields.
- A **subquery** may be useful to break your question down into subparts. Consider the following query skeleton, which uses the `WITH` operator to store a subquery's results in a temporary table named `donations`.

```
WITH donations AS (  
    SELECT ...  
    FROM ...  
    ... JOIN ...  
    ON ...  
    WHERE ...  
)  
SELECT ...  
FROM donations  
GROUP BY ...  
ORDER BY ...;
```

```
[ ]: # %%sql  
# /*  
#  * Code in this scratchwork cell is __not graded.__  
#  * Copy over any SQL queries you write here into the below Python cell.  
#  * Do __not__ insert any new cells in between the SQL/Python cells!  
#  * Doing so may break the autograder.  
#  */  
# -- Write below this comment. --
```

```
[ ]: query_q6 = """  
...     # replace this with  
...;     # your multi-line solution  
"""
```

```
res_q6 = pd.read_sql(query_q6, engine)
res_q6
```

```
[ ]: grader.check("q6")
```

## 5 Congratulations! You finished the lab!

### 5.1 Submission

Make sure you have run all cells in your notebook in order before running the cell below, so that all images/graphs appear in the output. The cell below will generate a zip file for you to submit. **Please save before exporting!**

```
[56]: # Save your notebook first, then run this cell to export your submission.
grader.export(pdf=False, run_tests=True)
```

Running your submission against local test cases...

Your submission received the following results when run against available test cases:

```
q1 results: All test cases passed!

q2 results: All test cases passed!

q3 results: All test cases passed!

q4 results: All test cases passed!

q5a results: All test cases passed!

q5b results: All test cases passed!

q6 results:
  q6 - 1 result:
    Trying:
      res_q6.shape == (10, 4)
    Expecting:
      True
*****
    Line 1, in q6 0
    Failed example:
      res_q6.shape == (10, 4)
    Exception raised:
      Traceback (most recent call last):
        File "/srv/conda/envs/notebook/lib/python3.9/doctest.py", line
1336, in __run
```

```

        exec(compile(example.source, filename, "single",
            File "<doctest q6 0[0]>", line 1, in <module>
            res_q6.shape == (10, 4)
        NameError: name 'res_q6' is not defined

q6 - 2 result:
    Trying:
        sorted(res_q6['total_amount']) == [1000, 2000, 2600, 5000, 5000,
5200, 5400, 9000, 10000, 18633157]
    Expecting:
        True
*****
        Line 1, in q6 1
        Failed example:
            sorted(res_q6['total_amount']) == [1000, 2000, 2600, 5000, 5000,
5200, 5400, 9000, 10000, 18633157]
        Exception raised:
            Traceback (most recent call last):
              File "/srv/conda/envs/notebook/lib/python3.9/doctest.py", line
1336, in __run
                exec(compile(example.source, filename, "single",
                File "<doctest q6 1[0]>", line 1, in <module>
                sorted(res_q6['total_amount']) == [1000, 2000, 2600, 5000,
5000, 5200, 5400, 9000, 10000, 18633157]
                NameError: name 'res_q6' is not defined

q6 - 3 result:
    Trying:
        sorted(res_q6['num_donations']) == [1, 1, 1, 1, 1, 1, 1, 2, 2,
131]
    Expecting:
        True
*****
        Line 1, in q6 2
        Failed example:
            sorted(res_q6['num_donations']) == [1, 1, 1, 1, 1, 1, 1, 2, 2,
131]
        Exception raised:
            Traceback (most recent call last):
              File "/srv/conda/envs/notebook/lib/python3.9/doctest.py", line
1336, in __run
                exec(compile(example.source, filename, "single",
                File "<doctest q6 2[0]>", line 1, in <module>
                sorted(res_q6['num_donations']) == [1, 1, 1, 1, 1, 1, 1, 2,
2, 131]
                NameError: name 'res_q6' is not defined

q6 - 4 result:

```

Trying:

```
sorted(res_q6['cmte_nm']) == [  
    'DONALD J. TRUMP FOR PRESIDENT, INC.',  
    'DONOVAN FOR CONGRESS',  
    'FRIENDS OF DAVE BRAT INC.',  
    'GRASSLEY COMMITTEE INC',  
    'HELLER FOR SENATE',  
    'NEW HAMPSHIRE REPUBLICAN STATE COMMITTEE',  
    'NY REPUBLICAN FEDERAL CAMPAIGN COMMITTEE',  
    'REPUBLICAN PARTY OF IOWA',  
    'SOUTH CAROLINA REPUBLICAN PARTY',  
    'TEXANS FOR SENATOR JOHN CORNYN INC']
```

Expecting:

True

\*\*\*\*\*

Line 1, in q6 3

Failed example:

```
sorted(res_q6['cmte_nm']) == [  
    'DONALD J. TRUMP FOR PRESIDENT, INC.',  
    'DONOVAN FOR CONGRESS',  
    'FRIENDS OF DAVE BRAT INC.',  
    'GRASSLEY COMMITTEE INC',  
    'HELLER FOR SENATE',  
    'NEW HAMPSHIRE REPUBLICAN STATE COMMITTEE',  
    'NY REPUBLICAN FEDERAL CAMPAIGN COMMITTEE',  
    'REPUBLICAN PARTY OF IOWA',  
    'SOUTH CAROLINA REPUBLICAN PARTY',  
    'TEXANS FOR SENATOR JOHN CORNYN INC']
```

Exception raised:

Traceback (most recent call last):

File "/srv/conda/envs/notebook/lib/python3.9/doctest.py", line 1336, in \_\_run

```
    exec(compile(example.source, filename, "single",  
File "<doctest q6 3[0]>", line 1, in <module>  
    sorted(res_q6['cmte_nm']) == [  
NameError: name 'res_q6' is not defined
```

<IPython.core.display.HTML object>