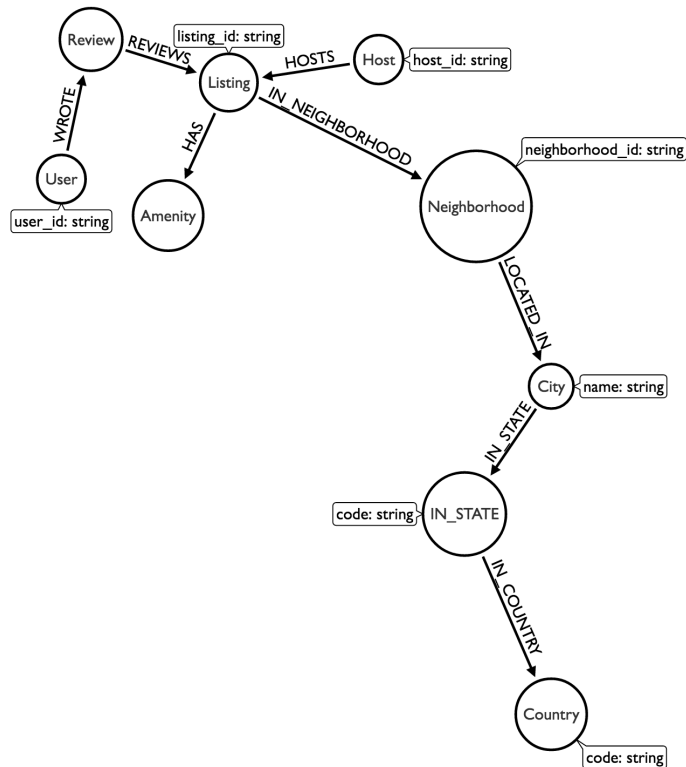


CS 327E Project 7, due Thursday 11/04.

This project makes use of an Airbnb dataset which is modeled as follows:



Your tasks are to load the Airbnb graph into your Neo4j database and then construct some queries on the resulting graph.

To start, open a new terminal window in JupyterLab and download and extract the airbnb assets from GCS:

```
gsutil cp gs://cs327e-open-access/airbnb.zip .
unzip airbnb.zip
```

The extracted folder contains 3 files: `listings.csv`, `reviews.csv`, `data_load.cypher`.

Second, create a new Python Jupyter notebook and name it `project6.ipynb`. All subsequent instructions should be run through your notebook unless otherwise noted.

Before loading any data, be sure to empty your Neo4j database by running this command:

```
!$CONNECT "MATCH (n) DETACH DELETE n"
```

Of course, you'll need to set the `CONNECT` variable before running the above command!

Third, load the airbnb data into Neo4j as follows:

```
!cat /home/jupyter/airbnb/load_data.cypher | {CONNECT} --format plain
```

The script may take a few minutes to run. It will print the number of nodes that it loads into Neo4j as it progresses.

Verify that all of the data has loaded correctly by returning a total node count:

```
!{CONNECT} "MATCH (n) RETURN count(n)"
```

You should get back 129,444 nodes.

If you have an SSH tunnel, you can bring up the [Neo4j Browser](#) and explore the nodes in the graph and their relationships. This step is not required. If you didn't bring up the Neo4j browser, you can refer to the Airbnb diagram above to see which node labels are in the graph and how they are connected.

Go back to your notebook and get a count of each unique node label in the graph:

```
!{CONNECT} "MATCH (n) RETURN distinct labels(n), count(n)"
```

You are now ready to construct some cypher queries. Start by sampling the data by returning the contents of 10 random nodes in one query and 10 random relationships in another query.

Next, translate these questions into cypher and output the results for each one.

- Q1. How many hosts are located in "Austin, Texas, United States"?
- Q2. Which listings does `host_id = "4641823"` have? Return the listing name, `property_type`, `price`, and `availability_365` sorted by price. Limit the results to 10.
- Q3. Which users wrote a review for `listing_id = "5293632"`? Return the user's id and name sorted alphabetically by name. Limit the results to 10.
- Q4. Which users wrote a review for any listing which has the amenities "Washer" and "Dryer"? Return the user's id and name sorted alphabetically by name. Limit the results to 10.

- Q5. Which listings have 3 bedrooms and are located in the Clarksville neighborhood? Return the listing name, property\_type, price, and availability\_365 sorted by price. Limit the results to 5.
- Q6. Which amenities are the most common? Return the name of the amenity and its frequency. Sort the results by count in descending order. Limit the results to 5.
- Q7. Which neighborhoods have the highest number of listings? Return the neighborhood's name and zip code (neighborhood\_id) along with the number of listings they have sorted by the number of listings in descending order. Limit the results to 5.

CS 327E Project 7 Rubric

**Due Date: 11/04/21**

Download and extract the airbnb assets to your jupyter notebook instance. -1 no airbnb assets found in Jupyter instance	1
Create a new Python Jupyter notebook named <code>project7.ipynb</code> . -1 incorrect file name	1
Run the data loader script ( <code>load_data.cypher</code> ) to populate the airbnb graph. -2 script not run or run incorrectly	2
Run a query that returns a count for each node label. -2 for missing or incorrect counts	2
Run a query that returns 10 random nodes and another that returns 10 random relationships. -3 for each missing or incorrect queries -2 for each missing or incorrect output from queries	10
Implement queries Q1 - Q7. -8 for each missing or incorrect query -4 for each missing or incorrect output from query	84
<code>project6.ipynb</code> pushed to your group's private repo on GitHub. Your project <b>will not</b> be graded without this submission.	<b>Required</b>
<code>submission.json</code> submitted into Canvas. Your project <b>will not</b> be graded without this submission. The file should have the following schema:  <pre>{   "commit-id": "your most recent commit ID from GitHub",   "project-id": "your project ID from GCP" }</pre> Example:  <pre>{   "commit-id": "dab96492ac7d906368ac9c7a17cb0dbd670923d9",   "project-id": "some-project-id" }</pre>	<b>Required</b>
<b>Total Credit:</b>	<b>100</b>