

Lab 2: Bacon Number

The questions below are due on Monday September 18, 2017; 10:00:00 PM.

You are not logged in.

If you are a current student, please Log In (<https://6009.csail.mit.edu/fall17/labs/lab2?loginaction=redirect>) for full access to the web site.

Note that this link will take you to an external site (<https://oidc.mit.edu>) to authenticate, and then you will be redirected back to this page.

Table of Contents

- 1) Preparation (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_1)
- 2) Introduction (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_2)
 - 2.1) The Film Database (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_2_1)
 - 2.2) The Names Database (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_2_2)
 - 2.3) Using the UI (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_2_3)
 - 2.4) `lab.py` and `test.py` (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_2_4)
- 3) Acting Together (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_3)
- 4) Bacon Number (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_4)
- 5) Paths (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_5)
 - 5.1) Bacon Paths (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_5_1)
 - 5.1.1) Speed (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_5_1_1)
 - 5.2) Arbitrary Paths (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_5_2)
- 6) Movie Paths (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_6)
- 7) Code Submission (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_7)
- 8) Feedback (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_8)
- 9) Checkoff (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_9)
 - 9.1) Grade (https://6009.csail.mit.edu/fall17/labs/lab2#catsoop_section_9_1)

1) Preparation

This lab assumes you have Python 3.5 or later installed on your machine.

The following file contains code and other resources as a starting point for this lab: `lab2.zip`

(https://6009.csail.mit.edu/fall17/lab_distribution.zip?

`path=%5B%22fall17%22%2C+%22labs%22%2C+%22lab2%22%5D`)

Most of your changes should be made to `lab.py`, which you will submit at the end of this lab. Importantly, you should not add any imports to the file.

This lab is worth a total of 4 points. Your score for the lab is based on:

- correctly answering the questions throughout this page (1 points)
- passing the `test.py` efficiently (2 points, see below), and
- a brief "checkoff" conversation with a staff member to discuss your code (1 points).

Your points for `test.py` are based on how quickly your code runs on the server:

- all tests correct and complete in less than 3.5 seconds each: 2 points
- all tests correct and complete in less than 9 seconds each: 1.5 points
- all tests correct and complete in less than 15 seconds each: 0.5 points

Note that each of the tests will be run in its own process.

2) Introduction

Have you heard of *Six Degrees of Separation*? This simple theory states that at most 6 people separate you from any other person in the world. (Facebook actually showed that the number among their users is significantly lower, at about 4.6. You can read more about it in a research paper (<http://arxiv.org/abs/1111.4570>).)

Hollywood has its own version: Kevin Bacon is the center of the universe (not really, but let's let him feel good about himself). Every actor who has acted with Kevin Bacon in a movie is assigned a "Bacon number" of 1, every actor who acted with someone who acted with Kevin Bacon is given a "Bacon number" of 2, and so on. (What Bacon number does Kevin Bacon have? Think about it for a second.)

Note that if George Clooney acts in a movie with Julia Roberts, who has acted with Kevin Bacon in a different film, George has a Bacon number of 2 through this relationship. If George himself has also acted in a movie with Kevin, however, then his Bacon number is 1, and the connection through Julia is irrelevant. We define the notion of a "Bacon number" to be the *smallest* number of films separating a given actor (or actress) from Kevin Bacon.

In this lab, we will explore the notion of the Bacon number. We have prepared an ambitious database of approximately 37,000 actors and 10,000 films so that you may look up your favorites. Did Julia Roberts and Kevin Bacon act in the same movie? And what does Robert De Niro have to do with Frozen? Let's find out!

2.1) The Film Database

We've mined a large database of actors and films from IMDB (<https://www.imdb.com/>) via the www.themoviedb.org (<http://www.themoviedb.org/>) API. We present this data set to you as a list of records (3-element lists), each of the form `[actor_id_1, actor_id_2, film_id]`, which tells us that `actor_id_2` acted with `actor_id_1` in a film denoted by `film_id`.

Keep in mind that "acts with" is a symmetric relationship. So if `[a1, a2, f]` is in the list, it is true both that `a1` acted with `a2` *and* that `a2` acted with `a1`, even if `[a2, a1, f]` is not explicitly represented in the list.

We store this data as JSON files (<https://en.wikipedia.org/wiki/JSON>). The server tests will use `small.json` and `large.json`, but we have also included a `tiny.json` that you will use to write your own tests.

2.2) The Names Database

The methods in `lab.py` expect you to use integer actor IDs, but the tests we give you on this page will have actor names as inputs and outputs.

We include a file, `resources/names.json`, which has a JSON (<https://en.wikipedia.org/wiki/JSON>) representation of the mapping between actor IDs and names. You can use the `load` method of Python's `json` module to get the data out of the file and into Python. For an example of this, check out how we load databases in the `setUp` functions of `test.py`.

Which of the following best describes the Python object that results from loading `resources/names.json`?

-- ▾

What is Terry Notary's ID number?

Which actor has the ID 1250109?

2.3) Using the UI

Your code will be loaded into a small server (`server.py`) and will serve up a visualization website. To use the visualization, run `server.py` and use your web browser navigate to `localhost:8000` (`http://localhost:8000/`). You will need to restart `server.py` in order to reload your code if you make changes.

You will be able to see actors as circular nodes (hover above the node to see the actor's name) and the movies as edges linking nodes together.

Above the graph, we define three different tabs, one for each component of the lab. Each tab sets up the visualization appropriate for its aspect of the lab.

2.4) `lab.py` and `test.py`

These files are yours to edit in order to complete this lab. You should implement the main functionality of the lab in `lab.py`, and you should implement additional test cases (as described throughout the assignment) in `test.py`.

In `lab.py`, you will find a skeleton for the functions we expect you to write.

3) Acting Together

Complete the definition of `did_x_and_y_act_together` in `lab.py`. This function should take three arguments, in order:

- The database to be used (a list of records of actors who have acted together in a film, as well as a film ID: `[actor_id_1, actor_id_2, film_id]`),
- Two IDs representing actors

This function should return `True` if the two given actors ever acted together in a film, and `False` otherwise. For example, Kevin Bacon (`id=4724`) and Steve Park (`id=4025`) did *not* act in a film together, meaning `did_x_and_y_act_together(..., 4724, 4025)` should return `False`.

Inside `test.py`, we have included a `TestTiny` class which has a `setUp` method but no tests. Add at least one test testing `did_x_and_y_act_together` on the tiny database (you can open `tiny.json` in a text editor to see the data it contains).

When you are done implementing this method and it passes the associated tests, use your code to answer the following questions according to the data in the `resources/small.json` database. (Hint: You will need to load `names.json` to get the mapping from actors to IDs, then find the actor corresponding to the ID. Consider writing a helper method that does this automatically.)

According to the `small.json` database, have Christian Stoeher and Craig Bierko acted together?

According to the `small.json` database, have Scott Subiono and Paule Annen acted together?

According to the `small.json` database, have Stephen Blackehart and Lew Knopp acted together?

According to the `small.json` database, have Philip Bosco and Patrick Malahide acted together?

Please note that `did_x_and_y_act_together` is a warmup and was given to help you to get familiar with the structure of the databases. You don't have to use the function in subsequent sections.

4) Bacon Number

Complete the definition of `get_actors_with_bacon_number` in `lab.py`. This function should take two arguments, in order:

- The database to be used (the same structure as before)
- The desired Bacon number

This function should return a Python list containing the ID numbers of all the actors with that Bacon number. The order of the elements doesn't matter. Note that we'll define the *Bacon number* to be the **smallest** number of films separating a given actor from Kevin Bacon, whose actor ID is `4724`.

Look at the data in `tiny.json` and answer these questions:

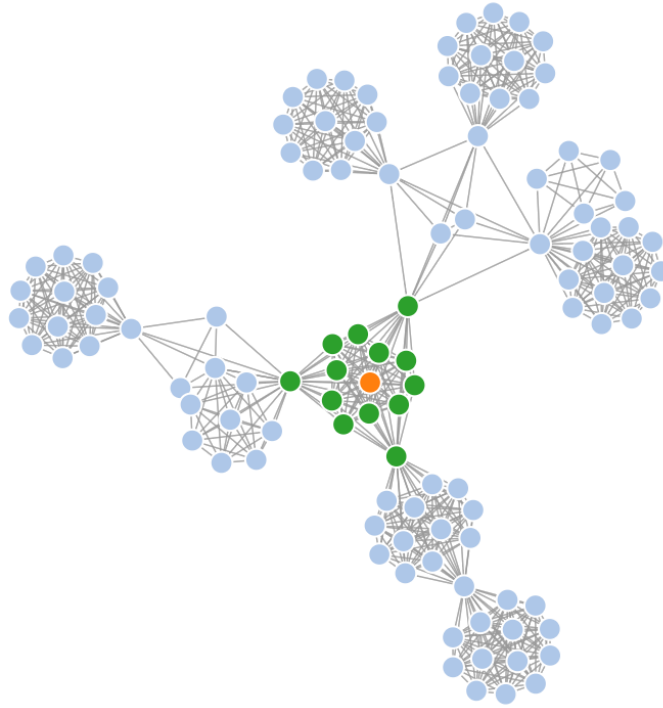
What are the IDs of the actors who have a Bacon number of 0 in `tiny.json`? Enter your answer below, as a Python list of integers:

What are the IDs of the actors who have a Bacon number of 1 in `tiny.json`? Enter your answer below, as a Python list of integers:

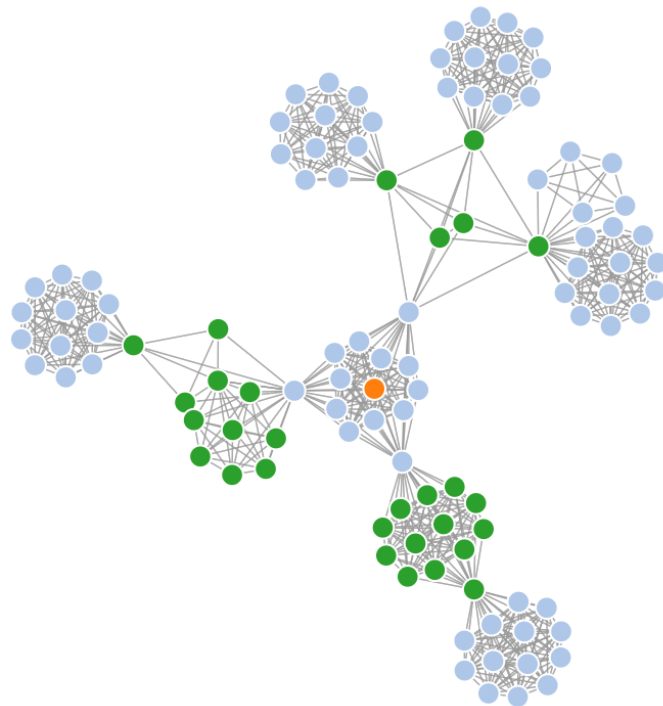
What are the IDs of the actors who have a Bacon number of 2 in `tiny.json`? Enter your answer below, as a Python list of integers:

Add the three above questions as tests in the `TestTiny` class, using your `get_actors_with_bacon_number` function to compute the lists of actors with Bacon numbers 0, 1, and 2; and comparing against the results you just found above.

Now you're ready to write your Bacon number code! Here are some things to think about when writing your implementation. Consider the set of actors with a *Bacon number* of 1. Here is a visual representation of the data from the `small.json` database:



Given the set of actors with a *Bacon number* of 1, think of how you can find the set of actors with a *Bacon number* of 2:



Once you get a sense for how to get the *Bacon number* 2 actors from the *Bacon number* 1 actors, try to generalize to getting the *Bacon number* $i+1$ actors from the *Bacon number* i actors.

Make sure that your function doesn't assign multiple Bacon numbers to the same actor.

Note that the test cases in `test.py` run against small and large databases of actors and films, and that your implementation needs to be efficient enough to handle the large database in a timely manner.

When you're done writing this method and it passed all of your tests, answer the following questions:

In the `small.json` database, what is the list of actors with Bacon number 3? Enter your answer below, as a Python list of *actor names*:

In the `small.json` database, what is the list of actors with Bacon number 4? Enter your answer below, as a Python list of *actor names*:

In the `large.json` database, what is the list of actors with Bacon number 5? Enter your answer below, as a Python list of *actor names*:

In the `large.json` database, what is the list of actors with Bacon number 6? Enter your answer below, as a Python list of *actor names*:

5) Paths

Now we'll turn our attention to finding the *chain* of actors that connects Kevin Bacon to someone else.

5.1) Bacon Paths

Complete the definition of `get_bacon_path` in `lab.py`. The function should take two arguments, in order:

- The database to be used (a list of records of actors who have acted together in a film, as well as a film ID: `[actor_id_1, actor_id_2, film_id]`),
- An ID representing an actor

Your function should produce a list of actor IDs (any such shortest list, if there are several) detailing a "Bacon path" from Kevin Bacon to the actor denoted by `actor_id`. If no path exists, return `None`.

Please note that the paths are not necessarily unique, and so any shortest list that connects Bacon to the actor denoted by `actor_id` is valid. The tester does not hard-code the correct paths and only verifies the *length* of the path you find (also that it is indeed a path that exists in the database).

For example, if we run this method with Julia Roberts's ID (`actor_id=1204`), one valid path is `[4724, 3087, 1204]`, showing that Kevin Bacon (`4724`) has acted with Robert Duvall (`3087`), who in turn acted with Julia Roberts (`1204`).

Take look at the data in `tiny.json`. What's the shortest path that connects actor `4724` to actor `46866`? Enter your answer below, as a Python list of ID numbers:

Add a test for the case above to the `TestTiny` class. You can use this test to help make sure your function is implemented correctly.

5.1.1) Speed

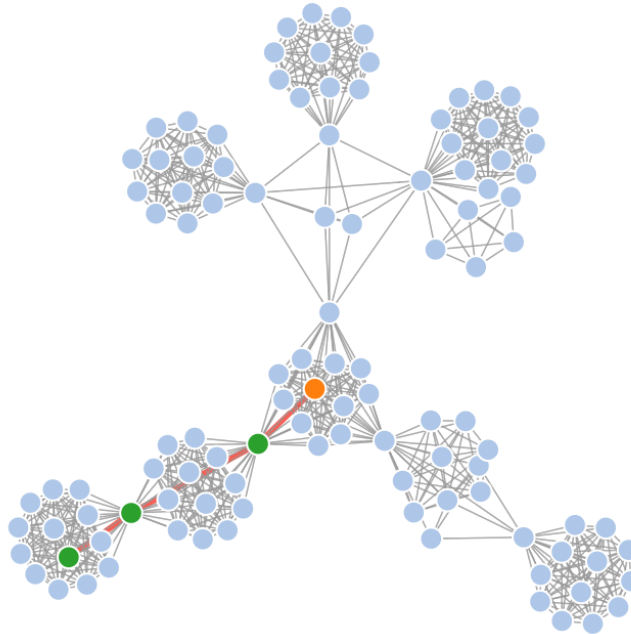
When implementing the path-finding algorithm, you should optimize your code to handle the large database, which our testing infrastructure will use when testing your code.

In particular, here are a few ideas about speed:

- Membership tests (the `in` operator) on long lists can be very slow. By contrast, the `in` operator is very fast on sets and dictionaries (regardless of the lengths of these objects). However, sets and dictionaries do not retain information about the order of their elements. Consider whether there are cases in your code where a set or dictionary can be used in place of a list.
- Running `L.pop(0)` on a long list is also slow. If you find yourself doing this, ask: do you really need to `pop`? Or can you just use an index to keep track of which list element you're working on?

- Searching through data using a `for` loop can be slow. Can you reorganize the data so that your search can be implemented with a single dictionary lookup or set containment check?

You will also need to be careful about your overall algorithm. In particular, **avoid repeatedly iterating through all of data**. For example, consider the following graph, with a path highlighted:



Here we've started from Kevin Bacon and successfully expanded out our search until we got to the actor we were looking for. What do we need to keep track of during our search if we want to get the path without looking for the actor again?

When you have implemented your function and it passes your tests, use it to answer the questions below:

According to the `large.json` database, what is the path of actors from Kevin Bacon to Ron Howard? Enter your answer as a Python list of actor names below:

According to the `large.json` database, what is the path of actors from Kevin Bacon to Jack Hoxie? Enter your answer as a Python list of actor names below:

According to the `large.json` database, what is the path of actors from Kevin Bacon to Anton Radacic? Enter your answer as a Python list of actor names below:

5.2) Arbitrary Paths

What we've done so far is pretty good, but it raises an important question: what makes Kevin Bacon so special? So far, everything we've done has centered around him... but let's expand things a bit, to be able to find the path that connects two *arbitrary* actors to each other.

Complete the definition of `get_path` in `lab.py`. The function should take three arguments, in order:

- The database to be used (a list of records of actors who have acted together in a film, as well as a film ID: `[actor_id_1, actor_id_2, film_id]`),
- Two IDs representing actors

Your function should produce a list of actor IDs (any such shortest list, if there are several) detailing a path from one actor to the other.

Add at least one test case for `get_path` to `TestTiny`, based on the contents of the `tiny.json` database. It should find the minimal path between two non-Bacon actors.

When you have implemented this function and it passes your tests, use it to answer the questions below:

According to the `large.json` database, what is the minimal path of actors from Beau Bridges to Vjeran Tin Turk? Enter your answer as a Python list of actor names below:

According to the `large.json` database, what is the minimal path of actors from Martin Ljung to Michael Yarmush? Enter your answer as a Python list of actor names below:

6) Movie Paths

After completing the work above, you might be interested to know what sequence of movies you could watch in order to traverse the path from one actor to another. For example, to move from Kevin Bacon to Julia Roberts, one could watch movie ID 94671 ("Jayne Mansfield's Car," which connects Kevin Bacon to Robert Duvall) and 18402 ("Something to Talk About," which connects Robert Duvall to Julia Roberts).

Add some code to your `lab.py` to determine the list of *movie names* that connect two arbitrary actors. To this end, we have included the `movies.json` database, which maps movie names to ID numbers.

When you have finished this code, use it to answer the following questions:

According to the `large.json` database, what is the minimal path of *movie titles* connecting Iva Ilakovac to Unknown Actor 15? Enter your answer as a Python list of movie names below:

According to the `large.json` database, what is the minimal path of *movie titles* connecting Jessica James to Ric Walker? Enter your answer as a Python list of movie names below:

7) Code Submission

When you have tested your code sufficiently on your own machine, submit your modified `lab.py` by clicking on "Choose File", then clicking the `Submit` button to send the file to the 6.009 server. The server will run the tests and report back the results (including timing) below.

Select File No file selected

8) Feedback

When you are finished with the lab, please help us out by answering the following questions:

How many hours, approximately, did you spend working on this lab?

On a scale of 1-7, how difficult was this lab?

-- ▼

Once you are finished with the code, please come to a tutorial, lab session, or office hour and add yourself to the queue asking for a checkoff. **You must be ready to discuss your code and test cases in detail before asking for a checkoff.**

- Your additional test cases in the `TestTiny` class.
- Your implementation of `did_x_and_y_act_together`.
- Your implementation of `get_actors_with_bacon_number`.
- Your implementation of `get_path` and `get_bacon_path`, and your test cases for `get_path`.
- How you transformed actor/movie names into ID numbers, and *vice versa*.
- Any additional code you wrote to compute the paths of actor/movie names.

You have not yet received this checkoff. When you have completed this checkoff, you will see a grade here.

Powered by CAT-SOOP (<https://catsoop.mit.edu/>) (development version).
CAT-SOOP is free/libre software (<http://www.fsf.org/about/what-is-free-software>), available under the terms of the GNU Affero General Public License, version 3 (https://6009.csail.mit.edu/cs_util/license).
(Download Source Code (https://6009.csail.mit.edu/cs_util/source.zip?course=fall17))