

Lab 2: Bacon Number

The questions below are due on Friday February 21, 2020; 04:00:00 PM.

You are not logged in.

If you are a current student, please [Log In](#) for full access to the web site.

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

Table of Contents

- [1\) Preparation](#)
- [2\) Introduction](#)
 - [2.1\) The Film Database](#)
 - [2.2\) The Names Database](#)
 - [2.3\) Using the UI](#)
 - [2.4\) lab.py and test.py](#)
- [3\) Acting Together](#)
- [4\) Bacon Number](#)
- [5\) Paths](#)
 - [5.1\) Bacon Paths](#)
 - [5.1.1\) Speed](#)
 - [5.2\) Arbitrary Paths](#)
- [6\) Movie Paths](#)
- [7\) Generalizing Our Path Finder](#)
 - [7.1\) Movie-to-Movie Paths](#)
- [8\) Code Submission](#)
- [9\) Checkoff](#)
 - [9.1\) Grade](#)

1) Preparation [§](#)

This lab assumes you have Python 3.6 or later installed on your machine.

The following file contains code and other resources as a starting point for this lab: [lab2.zip](#)

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.

You can also see and participate in online discussion about this lab in the "[Lab 2](#)" [Category](#) in the forum.

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 point)
- passing the `test.py` tests (2 points, see below), and
- a brief "checkoff" conversation with a staff member to discuss your work, and a review of your code's clarity/style (1 point).

Your points for the test cases are based on how quickly your code runs on the server together with correctness. Your lab code will be expected to run in a certain amount of time (which varies depending on the test); this execution time is based on the amount of time it takes to run that test on the server (*not* your own machine).

The questions in sections 1-3 (inclusive) on this page are due before lecture, at 1:30pm on Tuesday, 18 Feb (which runs on a Monday schedule). The remaining questions on this page (including your code submission) are due at 4pm on Friday, 21 Feb. Checkoffs are due at 10pm on Wednesday, 26 Feb.

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.

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 API. We present this data set to you as a list of records (3-element tuples), 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. If `(a1, a2, f)` is in the database, 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 database.

However, these relationships do not necessarily exhibit the transitive property. That is, if `(a1, a2, f)` and `(a2, a3, f)` are in the database, it is *not necessarily true* that `a1` and `a3` have acted together (unless `(a1, a3, f)` or `(a3, a1, f)` is in the database).

We store these data as [pickle files](#). The server tests will use `small.pickle` and `large.pickle`, but we have also included a `tiny.pickle` 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.

To help with this mapping, we include a file, `resources/names.pickle`, which contains a representation of the mapping between actor IDs and names. You can use the `load` method of Python's `pickle` 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` function of `test.py`.

Answer the following questions *using Python*. Even though these are small snippets, please include them in the `if __name__ == '__main__':` block in your `lab.py` file so that we can discuss them during your checkoff.

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

This question is due on Tuesday February 18, 2020 at 01:30:00 PM.

What is Hugh Hurd's ID number?

This question is due on Tuesday February 18, 2020 at 01:30:00 PM.

Which actor has the ID 76355?

This question is due on Tuesday February 18, 2020 at 01:30:00 PM.

2.3) Using the UI

We have also provided a visualization website which loads your code into a small server (`server.py`) and visualizes your results. To use the visualization, run `python3 server.py` and use your web browser navigate to [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 skeletons for most of the functions we expect you to write.

3) Acting Together

To get used to the structure of the databases, complete the definition of `acted_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, of the form described above),
- 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 `acted_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 `acted_together` on the tiny database (you can load the data from `tiny.pickle` and print the results to see what 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.pickle` database. (Hint: You will also need to load `names.pickle`, to determine an actor's ID number from their name.)

According to the `small.pickle` database, have Joseph McKenna and Dan Warry-Smith acted together?

This question is due on Tuesday February 18, 2020 at 01:30:00 PM.

According to the `small.pickle` database, have Josef Sommer and Stig Olin acted together?

This question is due on Tuesday February 18, 2020 at 01:30:00 PM.

Please note that `acted_together` is intended to provide a way to help you get familiar with the structure of the databases. You don't have to use the function in subsequent sections. Also note, though, that `test.py` does test this function.

4) Bacon Number

Next, we will try to find all of the actors who have a given Bacon number. We'll implement this as a function called `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 set containing the ID numbers of all the actors with that Bacon number. 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.

Before we get to writing this function, we'll develop a couple of test cases that we can use to make sure our function works properly (once we've written it!). Look at the data in `tiny.pickle` and answer these questions (by computing the responses manually from looking at the data):

What are the **ID numbers** of the actors who have a Bacon number of 0 in `tiny.pickle`? Enter your answer below as a Python set of integers:

This question is due on Friday February 21, 2020 at 04:00:00 PM.

What are the **ID numbers** of the actors who have a Bacon number of 1 in `tiny.pickle`? Enter your answer below as a Python set of integers:

This question is due on Friday February 21, 2020 at 04:00:00 PM.

What are the **ID numbers** of the actors who have a Bacon number of 2 in `tiny.pickle`? Enter your answer below as a Python set of integers:

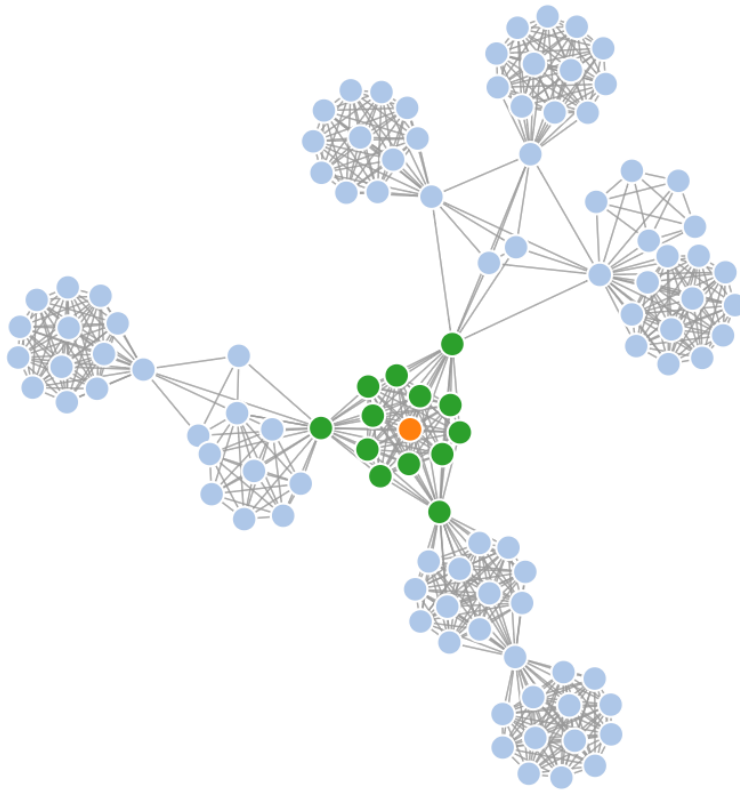
This question is due on Friday February 21, 2020 at 04:00:00 PM.

What are the **ID numbers** of the actors who have a Bacon number of 3 in `tiny.pickle`? Enter your answer below as a Python set of integers:

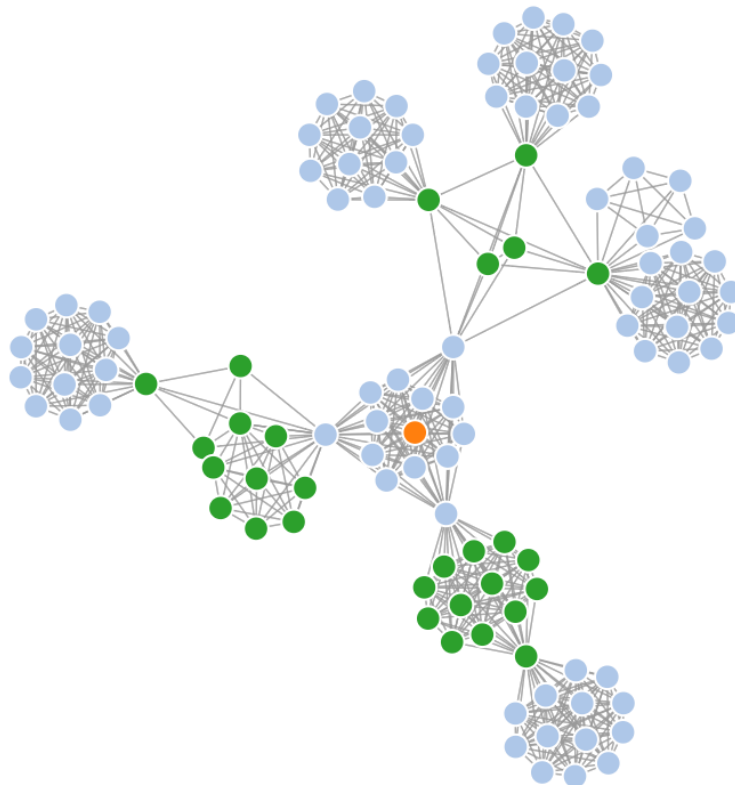
This question is due on Friday February 21, 2020 at 04:00:00 PM.

Add the four above questions as tests in the `TestTiny` class. In each test case, use your `actors_with_bacon_number` function to compute the sets of actors with Bacon numbers 0, 1, 2, and 3 and compare them against the results you just found above. So that they are recognized as individual tests, **define each as a separate method in the class, starting with `test_`** (e.g. `test_bacon_number_0`).

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.pickle` database. (You can use our provided server to generate pictures like these.)



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.

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. Your code should also handle the case of arbitrary Bacon numbers (not just $n \leq 6$) since some databases may be structured to assign some actors quite large Bacon numbers.

When you're done writing this method and it passes all of your tests, answer the following question that uses the `large.pickle` database:

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

This question is due on Friday February 21, 2020 at 04:00:00 PM.

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 `bacon_path` in `lab.py`. The function should take two arguments in order:

- The database to be used (the same structure as before),
- 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, 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 (as well as that it is indeed a path that exists in the database).

For example, if we run this method on `large.pickle` 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 a look at the data in `tiny.pickle` (by loading and printing it, not by running your path-finding code). From those data, you should be able to manually compute a couple of shortest paths. What's the shortest path that connects actor 4724 to actor 1640? Enter your answer below as a Python list of **ID numbers**:

This question is due on Friday February 21, 2020 at 04:00:00 PM.

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

implemented correctly.

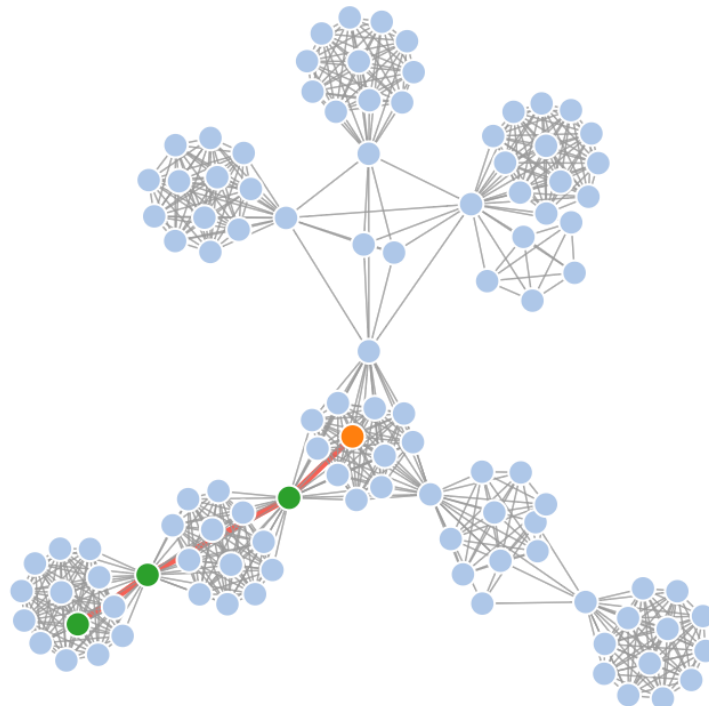
5.1.1) Speed

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:

- 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?
- 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?

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 question below:

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

This question is due on Friday February 21, 2020 at 04:00:00 PM.

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. Let's expand things a bit and find the path that connects two *arbitrary* actors to each other.

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

- The database to be used (the same structure as before),
- 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 the first actor to the second.

Add at least one test case for `actor_to_actor_path` to `TestTiny` based on the contents of the `tiny.pickle` 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 question below. Even though some of these may be small code snippets, please include them in the `if __name__ == '__main__':` block of your `lab.py` so that we can discuss them during your checkoff.

According to the `large.pickle` database, what is the minimal path of actors from Awie to Steve Carell? Enter your answer as a Python list of **actor names** below:

This question is due on Friday February 21, 2020 at 04:00:00 PM.

Now that you have implemented it, take a look at your definition of `actor_to_actor_path`. How does it compare to `bacon_path`? You may find that they are very similar; in that case, it may be a good idea to "refactor" your code (rearrange it, perhaps introducing helper functions) to avoid repetitious code where possible. As you are refactoring, it is a good idea to continue testing your functions after each big change, to make sure they are still working as expected.

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.pickle` database, which maps movie names to ID numbers. **You should be prepared to discuss this code during your checkoff.**

When you have finished this code, use it to answer the following question. Include your code for answering this question in your `lab.py` file so that we can discuss it during the checkoff.

According to the `large.pickle` database, what is the minimal path of *movie titles* connecting Tracy Reiner to Sven Batanic? Enter your answer as a Python list of **movie names** below:

This question is due on Friday February 21, 2020 at 04:00:00 PM.

7) Generalizing Our Path Finder

As it currently stands, our `actor_to_actor_path` function can currently only search for a particular actor. But it turns out that there are certain kinds of questions we could ask that are, unfortunately, difficult to ask given this implementation.

For example, suppose we want to find the shortest path from some actor to *any actor from a set of other actors*, or from some actor to *any actor in a particular movie*, or something like that.

To answer these kinds of questions, it may be helpful to generalize our notion of path-finding. Do so by filling in the definition of `actor_path` in `lab.py`. This function should take three arguments in order:

- The database to be used (the same structure as before),
- One actor ID to be used as our starting point, and
- A *function* to be used as our goal test. This function should take a single actor ID as input, and it should return `True` if that actor represents a valid ending location for the path, and `False` otherwise.

Your function should produce as output a list containing actor IDs, representing the shortest possible path from the given actor ID to *any actor that satisfies the goal test function*. If no actors satisfy the goal condition, your function should instead return `None`.

Note that if the starting actor satisfies the goal test, your function should return a length-1 list containing only that actor's ID.

7.1) Movie-to-Movie Paths

As our final task for this lab, we would like for you to find chains of actors that connect one given *movie* to another. Implement this behavior by filling in the body of the `actors_connecting_films` function. This function should take three arguments:

- The database to be used (the same structure as before), and
- Two film ID numbers;

and it should return the shortest possible list of actor ID numbers (in order) that connect those two films. Your list should begin with the ID number of an actor who was in the first film, and it should end with the ID number of an actor who was in the second film.

If there is no path connecting those two films, your function should return `None`.

8) Code Submission

When you have tested your code sufficiently on your own machine, submit your modified `lab.py` below. Note that your checkoff (including style considerations) will be based on your most recent submission, and that all aspects of the file will be considered in terms of style, including those that are not explicitly tested in `test.py` (for example, any helper functions you write).

When submitting `lab.py`, the server will run the tests and report back the results (including timing). Submit your `lab.py` in the box below:

No file selected

This question is due on Friday February 21, 2020 at 04:00:00 PM.

Submit your `test.py` below, and **be prepared to discuss it during the checkoff**. Your `test.py` will be run on a few different implementations of the lab (some of which are correct, and some of which are not).

No file selected

This question is due on Friday February 21, 2020 at 04:00:00 PM.

9) Checkoff

Once you are finished with the code, please come to a lab session or office hour and add yourself to the queue asking for a checkoff. **You must be ready to discuss your code in detail before asking for a checkoff.** Since the clarity of your code will be evaluated as part of the checkoff, you may wish to take some time to comment your code, use good variable names, avoid repetitive code (create helper methods), etc.

Be prepared to discuss:

- Your additional test cases in the `TestTiny` class.
- Your implementation of `acted_together`.
- Your implementation of `actors_with_bacon_number`.
- Your implementation of `actor_path`, `actor_to_actor_path`, and `bacon_path`, and your test cases for `actor_to_actor_path`.
- How you transformed actor/movie names into ID numbers, and *vice versa*.
- The additional code you wrote to compute the paths of actor/movie names.

9.1) Grade

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