

<http://www.cs.arizona.edu/classes/cs227/spring07/>

Program #10: Binary Searching a Binary File

Due Date: April 23rd, 2007, at 10:00 p.m. MST

Overview: Binary Search is normally presented as a technique for searching in-memory data structures. However, it can be effectively used to search for uniformly-sized records in binary files, too.

Thanks to your taxpayer dollars and the requirement that the government conduct a census of the population every 10 years, the U.S Census Bureau collects a large volume of data and makes much of it available to anyone who cares to view it. For this assignment, a bit of it makes great data for an ordered binary file of records. In the source folder on your X: drive you should find a binary file named `county2k.bin` with some info on over 3,200 counties from the 2000 census. Each 78-byte record of that file has data about a specific county (or equivalent unit) in the United States. In this assignment, you'll let the user search the `county2k.bin` file (using a recursive binary search) for county names.

Assignment: Write a complete, well-documented Java program named `Prog10.java` that asks the user for the complete path to the `county2k.bin` file and the name of a county, and, using a recursive binary search, finds (or fails to find) *all* counties within that file whose names begin with that name prefix (case doesn't matter; 'A' matches 'a'). If the name is found, the program will display to the screen the name, state abbreviation and population of *each* county starting with the given name prefix. If the name is not found, the program will display a message to that effect. Additional output requirements are detailed in the "Output" section, below.

Please be careful to note that the binary search is to be on the `county2k.bin` file of records. You may **NOT** write your program to read the records from the file into an array and search the array.

Data: The file `county2k.bin` can be found in the source folder on your X: drive, visible from the machines in the lab. It's also linked to the class web page. If you save it to another computer, make sure that your saved copy is the same size – exactly 251,082 bytes.

Each county and county-like unit is represented in the file by five fields of data totaling 78 bytes. In order, those fields are: A 2-byte state abbreviation (ex: AZ), an `int` value for the state FIPS code (ex: 4), an `int` value for the county FIPS code (ex: 19), an `int` value for the county population (ex: 843746), and a 64-byte string with the name of the county (ex: Pima County). Java's classes for reading from binary files provide methods for reading integers and characters.

You may be wondering how you can conduct a binary search without an array. We know that to conduct a binary search, we need data in sorted order and in a direct-access data structure. The binary file has been created with the county records in order by name, meeting the first condition. Because our county records are each 78 bytes in size, we can compute the first byte of the n -th record just as we computed the first byte of the n -th element of an array: The first record starts at byte 0, the second at byte 78, etc. All we need to do is jump to the computed byte, which Java provides methods to do (see the Hints section, below), followed by reads of the fields of the record that starts there.

Output: In addition to the list of county names, state abbreviations and populations (for a successful search) or the failure message (for an unsuccessful search), output the following. As the binary search is performed, for each 'probe' into the binary file, display to the screen the values of 'low,' 'high,' and 'mid' (think of these as file record numbers, where the first record in the file is 0), along with the county names of the records associated with each of the three. This will help you (and us!) see how the search progressed. Whether the binary search has succeeded or failed, display the quantity of 'mid' records that were read from the file. This count is not to include the reads of the 'low' and 'high' records, nor the additional reads necessary to find and display any additional counties with the same name.

(Continued ...)

Here's an output fragment. You don't have to follow this format, just display all of the requirement information in a clear manner:

```
Enter the binary file location: X:\Source\county2k.bin
Enter a county name prefix: barb

Probe #1:  Low:      0 (Abbeville County)
           Mid:    1609 (Lake County)
           High:   3218 (Ziebach County)

[...]

The counties whose names start with the letters of 'barb' are:

    KS      5307  Barber County
    AL     29038  Barbour County
    WV     15557  Barbour County

The number of 'mid' records read was [...]
```

Note that “Santa Barbara County” isn't included; its name contains ‘barb’, but doesn't begin with it.

Turn In: Electronically submit your `Prog10.java` file using the `turnin` program and the usual `[section leader first name]_S[your section #]_P[program #]_[your last name]_[your first name]` naming convention.

Want to Learn More?

- “FIPS” stands for **F**ederal **I**nformation **P**rocessing **S**tandard.
- www.census.gov is the U.S Census Bureau's home page.
- The `county2k.bin` file was created from the `county2k.txt` file available from the Gazetteer: <http://www.census.gov/geo/www/gazetteer/places2k.html>

Hints, Reminders, and Other Requirements:

- We're asking you to prompt for the binary file's location so that the SLs don't have to change your code to grade your program.
- A record in a file is just a logical concept describing a physical group of fields. In this assignment, a binary file record has five fields. By reading the five fields in order, you've read a record. That's all there is to it.
- It's been a while, so a reminder is in order: We covered binary file I/O earlier in the semester. Reviewing the example programs, particularly `T03n05.java` and `T03n06.java`, will be helpful.

You may wish to look at one class in particular that we didn't cover when we talked about binary files: `RandomAccessFile`. It's nice because you can use it to read from as well as write to binary files. Plus, it has a convenient `seek()` method. You don't need to use `RandomAccessFile`, but you are welcome to use it if you'd like to.

- We suggest that you start by writing a small program that reads the content of the first record and displays that information to the screen. This will demonstrate that you've figured out how to read a record. Then, extend that program to read all of the records. Next, try reading the records in reverse order (record n , then record $n - 1$, etc.). Only when you can do these things should you worry about implementing the recursive binary search.
- Your program will work if the number of counties changes, won't it?