

CSC 212—Algorithms and Complexity

Compare the Boyer-Moore (BM) and Knuth-Morris-Pratt (KMP)

pattern-matching methods.

Term 3 Practical 3

27 July 2015

Submit as soon as possible before leaving the practical.

Add the Boyer-Moore (BM) and Knuth-Morris-Pratt (KMP) pattern matching methods to last week's code.

Use the the same data as for last week's practical.

Copy your previous working practical in `23practical` into today's directory `33practical`. Add the Boyer-Moore (BM) and Knuth-Morris-Pratt (KMP) pattern matching methods. Use the large set of words in the file

`/export/home/notes/ds/lesmiserables.text`

in the SunLab as the file in which you must look up and count appearances of these words that we call patterns. A supply of such patterns to be used is stored in the file `patterns.text` in the same directory.

First the entire text file must be read into the array `String T` and then the patterns must be read *one-by-one* into the variable `String p`. Each pattern must be looked up using one of each of the methods. For each method the program must run for about 1 second and then report how many patterns were found in the time of one second. Your pattern matching routines must search repeatedly for a pattern—a number—until all its occurrences have been found and report on the number of its appearances. Write this to a file.

The final output of your code to the screen is a table that gives the name of each method and the number of patterns whose first occurrence it has matched within 1 second.

We advise you to test your code by checking if it works for the first 10 (ten) patterns. When this works, proceed to do the timing. Do all the timing for one method before proceeding to the next one. If you read any ambiguity into this problem statement please resolve it yourself.

1. *Preparation of file system:* Today's practical must be done as usual inside your name file that lies in your home directory. Do today's work in a directory called `33practical`. Create the new directory by copying your `23practical` directory into the new one. Fix up the `Makefile` so that it processes today's work.
2. **Today:** First create a Boyer-Moore (BM) pattern searching method that works for a fixed piece of text called `T` and a smallish but fixed pattern called `p`.
3. **Today:** Next create a Knuth-Morris-Pratt (KMP) pattern searching method that works for a fixed piece of text called `T` and a smallish but fixed pattern called `p`.
4. **Today:** Finally get both methods for pattern searching to work with our data, test them with smallish data, and when the methods have been shown to work, run the timed experiment. The work you need to do is to figure out how the algorithms work and get them to work. This must be *written up* and the write up must be handed in on or before 3 August 2015. This must be done by opening a directory `33-writeup` in which your L^AT_EX documentation must be stored.

Your coding skills will improve if the code is always your own. Do your debugging by first getting each method to work with a small text `T` and pattern `p` that are assigned in

the test code before trying to apply the search methods on the larger text and patterns from the files. Also do tests that (1) succeed at the start of T, (2) somewhere away from the end points of T, and (3) one that matches at the end of T, for example:

```
String p = "The";
String T = "The Rabin-Karp method is much faster than the brute-force method.";

String p = "brute-force";
String T = "The Rabin-Karp method is much faster than the brute-force method.";

String p = "method.";
String T = "The Rabin-Karp method is much faster than the brute-force method.";

String p = "";
String T = "The Rabin-Karp method is much faster than the brute-force method.";

String p = "method.";
String T = "";

String p = "";
String T = "";
```

Please note that these examples are textual, but the actual practical uses strings that consist entirely of strings separated by spaces or newline characters. Once you are convinced that your code works for these examples the program must be enhanced to read a pattern `p` from the `patterns.text` file and search for that pattern once in the text `T` that has also been read from a file. *Do not copy the data.* Simply read it and use it to generate your results. The data is stored in the files given above.

- Time the BM method for 1 second and report the number n of patterns that were searched.
 - Do the same for the KMP method.
5. Submit ALL your pracs in a file that has been created in the usual manner: i.e.,
`cd; make submit.`
 6. **The pattern drawing strategy.** Draw the patterns from the `patterns.text` file. For each $i \in [1..P]$, draw a pattern p_i from the `patterns.text` file and then search through the text for its occurrence in the `lesmiserables.text` file. Note that some patterns might not be present in the text file.
 7. Note that all the occurrences of a pattern must be looked up and counted.
 8. In summary look up each “pattern” p throughout the entire “text” file T .
-