

## CSE 102 Homework Assignment 5 (Due: November 9, 11:55 pm)

You are going to write a complete C program which implements the following functionality:

- Your program will read the following files: language\_1.txt language\_2.txt language\_3.txt language\_4.txt language\_5.txt language\_x.txt
- Each file contains text in a specific language. All files contain only english lowercase characters and whitespace. Text files will include the following characters: 'a' 'b' 'c' 'd' 'e' 'f' 'g' 'h' 'i' 'j' 'k' 'l' 'm' 'n' 'o' 'p' 'q' 'r' 's' 't' 'u' 'v' 'w' 'x' 'y' 'z' ' '
- Your program will evaluate the dissimilarity scores of language pairs:

(language\_x, language\_1)

(language\_x, language\_2)

(language\_x, language\_3)

(language\_x, language\_4)

(language\_x, language\_5)

- First of all, calculate bi-gram frequencies for each language. A bi-gram is defined as follows: For a given sequence, each unique pairing of successive letters is a bi-gram. For example: for the word "adana " bi-grams are defined to be " a", "ad", "da", "an", "na", "a ". Beware: If there is a space before or after a character you will still be dealing with bi-grams. Each bi-gram has exactly two elements which are either characters or space. In order to calculate the frequency of a particular bi-gram (lets say bi-gram "ad") you have to count all the bi-grams in a given text and for this bi-gram calculate the ratio (# of "ad")/(total # of all bi-grams)
- Given all the frequencies, dissimilarity score is calculated as follows:  $\text{dissimilarity}(\text{language}_a, \text{language}_b) = \sum_i |f_{i,a} - f_{i,b}|$  (1)
- Here  $f_{i,a}$  represents the frequency of  $i$ th bi-gram for the language  $a$ . If  $c_{i,a}$  is the count of  $i$ th bi-gram in language  $a$ , then;

$$f_a^i = c_a^i / (\sum_j c_a^j)$$

- After evaluating dissimilarities, your program will print all the dissimilarity values. Print:  $\text{dissimilarity}(\text{language\_x}, \text{language\_1})$   $\text{dissimilarity}(\text{language\_x}, \text{language\_2})$   $\text{dissimilarity}(\text{language\_x}, \text{language\_3})$   $\text{dissimilarity}(\text{language\_x}, \text{language\_4})$   $\text{dissimilarity}(\text{language\_x}, \text{language\_5})$  Remarks:
- text files can include multiple concatenating whitespace. For example: Here we are using a user defined recursive
- Two adjacent whitespace do not create a bi-gram.

- Input files can be multi-line text files.
- There isn't any limit on the size of input files. Your program should work regardless of the size of the input. Hints:
- Bi-gram types do not depend on the input file. There are finite number of possibilities. Given all the lowercase english characters and a space, you can generate all the possible bi-grams.
- Do not try to store all the content of the file in the memory. Counting is possible without storing all of the text.
- You don't need to parse words. Turn in: A complete C program which can be compiled using the following command: `gcc -std=c99 assignment_5_name_id.c -o assignment_5_name_id` If your program requires additional compile and link options, state that requirement at beginning of your source code as a comment. 2 Caution:
- Read and apply "Assignment Submission Rules and Other Related Information" document which available on the class e-learning system.
- You may or may not get partial credit depending on how you structured or documented your code.