

T9 Predictions Based on Unigram Frequencies



On the old Nokia and Verizon phones, when the complete keyboard was absent, users were expected to use their numeric keypad with just 9 keys to type in all 26 characters of the English alphabet.

The keys and the corresponding letters were

```
2 abc
3 def
4 ghi
5 jkl
6 mno
7 pqrs
8 tuv
9 wxyz
```

If a person types in '2' he could mean either 'a' or 'b' or 'c'.

If a person types in 23 he could mean either [*ad* or *ae* or *af* or *bd* or *be* or *bf* or *cd* or *ce* or *cf*].

In this problem, we try to guess, that if the person keys in a particular sequence of numeric keys on the keypad, what is the most likely word which he was trying to enter?

An outline of steps for building a Unigram model

1. A dictionary of commonly used words be provided to you.
2. Apart from that you are also provided a large corpus of text. Using this corpus of text, you can compute the frequency with which commonly used words (from the dictionary) occur in the corpus. i.e., you are computing Unigram Frequencies, using the corpus. There might be words in the corpus which are missing in the dictionary: these words can be ignored. But do not ignore words which are present in the dictionary and absent in the corpus.

Let's say, the dictionary contains the words

```
Hello,I,am
```

and suppose the corpus text is

```
Hello, I will be going to the mall tomorrow. I am out of groceries.
```

The weighted dictionary (with frequencies) you will come up with is:

```
[Hello=>2,I=>3,am->2]
```

Because *Hello* occurs once in the corpus text and once in dictionary, 'I' occurs twice in the corpus and once in dictionary and 'am' occurs once in corpus and dictionary.

Predicting the word, given a series of numerals

After you read the dictionary and corpus and build the language model, you are given a number of numeric sequences typed in by a phone user. You need to identify words from the dictionary which start with a character sequence which could be represented by these numerals. Identify the top five candidates with the highest frequency, and output them in one line, separated by semi-colons. If there are less than five possible candidates, display them all. If there is no possible candidate, display

No Suggestions

Most likely word is the one which occurs max times in the given corpus, least likely is the one which occurs least times in the given corpus (or, perhaps it is a word which exists only in the dictionary and did not occur at all in the corpus).

Dictionary and Corpus File

For the purpose of building the word frequency and unigram model, you are provided with a file *t9Dictionary* and *t9TextCorpus* which will be kept in the same folder as the one from which your program is being run.

1. The first file *t9Dictionary*, is the dictionary. First line contains N, N words follow each in a new line.
2. The second file is the training corpus *t9TextCorpus* of text. This ends with "END-OF-CORPUS" on a new line.

Defining a word

A word is a sequence of characters (a-z, lowercase or uppercase; hyphen or apostrophe) which always starts and ends with a letter (a-z, lowercase or uppercase). The regex used must be greedy.

Input Format

First line will contain the number of tests T. This will be followed by T lines containing numeric sequences/numbers N.

Constraints

```
1 <= T <= 20
2 <= N <= 10^10
```

Each digit in N is between 1 to 9.

Output Format

Each line will contain a list of top five semi-colon separated candidate words, with the leftmost word being the most frequently used. If the group of five words contains multiple words with the same frequency, sort them in lexicographical order. If there were less than five candidates for some input, display all of them. If there were no candidates for some input sequence, display "No Suggestions".

Sample Input

```
6
6837
86
69
23777
11111
77777
```

Sample Output

```
over;overlying;overcome;overcoat;overgrowth
to;under;took;united;too
my;own;myself;owing;owners
cesspool;cesspool's;cesspools
No Suggestions
No Suggestions
```

Explanation

6837, means that we are looking for words where the first character is either 'm' or 'n' or 'o'

second character is either 't' or 'u' or 'v' third character is either 'd' or 'e' or 'f' fourth character is either 'p' or 'q' or 'r' or 's'.

We select only the words from the dictionary, which match the above criteria, and we sort them in descending order of their frequencies, and display the first five. 'over' is the most frequent one, followed by 'overlying' and so on.

Similarly, we process other inputs.

In the last case, we don't have any words in the dictionary corresponding to the numeric sequence 77777, so we display "No Suggestions".

Scoring

Your score for a test case will be

$(\text{maxScore for the test case}) * (\text{nCorrect}) / (\text{nTotal})$

nCorrect = Number of lines matching with the expected output

nTotal = Total number of lines in the expected output