

A Bangla Phonetic Encoding for Better Spelling Suggestions

Naushad UzZaman and Mumit Khan
BRAC University, Dhaka, Bangladesh

naushad@bracuniversity.ac.bd, mumit@bracuniversity.ac.bd

Abstract

We present a phonetic encoding for Bangla that can be used by spelling checkers to provide better suggestions for misspelled words. The encoding is based on the Soundex algorithm, modified to match Bangla phonetics. We start by analyzing Soundex encoding scheme when applied to Bangla. Next we propose a new encoding that handles the case of Bangla words, including those containing conjuncts. We conclude with a demonstration of a prototype spelling checker that uses this phonetic encoding to offer suggestions for a set of misspelled Bangla words.

1. Introduction

One of the more difficult tasks for a spelling checker is to produce “good” suggestions for misspelled words. While there have been significant research efforts in approximate string matching algorithms for English and other Western languages [1-4], similar work for Bangla has however just begun [5, 6]. An analysis of Bangla misspelled words shows that two of most common reasons for misspellings are (i) phonetic similarity of Bangla characters, and (ii) the difference between the grapheme representation and phonetic utterances [7]. This observation is the primary motivation for creating a phonetic encoding for Bangla that can be used to provide suggestions for misspelled words. While this paper focuses on the spelling checking application, the proposed encoding is equally applicable in a wide range of text-processing applications, from searching for patient records in a medical database to matching names in census records. The basic idea behind spelling suggestions using phonetic encoding is quite simple:

1. Encode the input word using phonetic coding rules;
2. Look up a phonetically encoded lexicon for words with the same code; and
3. Create an ordered list, i.e., suggestions, from the result using some heuristic.

In this paper, we introduce a phonetic encoding for Bangla, and then demonstrate how a spelling checker would use it to produce suggestions for

misspelled Bangla words. We assume that the Bangla text is encoded using Unicode Normalization Form C (NFC) [9], with its consistent logical ordering of the consonants and the dependent vowels, as well as of the large repertoire of the *juktakkhors* (compound letters or conjuncts) in Bangla.

2. Phonetic matching techniques

A major class of approximate string matching algorithms are the various phonetic methods, from the eighty-year old Soundex [9, 10], to the more recent Metaphone [11, 12] and PHONIX [13]. The input to these phonetic encodings or “sound-alike” algorithms is a word, and the result is an encoded key, which should be the same for all words that are pronounced similarly, allowing for a reasonable amount of fuzziness. The basic principle behind these phonetic matching schemes is to partition the consonants by phonetic similarity, and then use a single key to encode each of these sets. Strings that sound similar compare equal in their respective encoded form. For these particular algorithms, only the first few consonant sounds are encoded, unless the first letter is a vowel. Metaphone for example encodes “Stephan”, “Steven”, and “Stefan” as STFNN, so all three names compare equal when encoded.

Of these phonetic methods, Soundex method is by far the oldest, first patented by Odell and Russel in 1918. Soundex partitions the set of letters into seven disjoint sets, assuming that the letters in the same set have similar sound. Each of these sets is given a unique key, except for the set containing the vowels and the letters h, w, and y, which is considered to be silent and is not considered during encoding. The Soundex codes are shown in Table 1. The Soundex algorithm itself, shown in Figure 1, transforms all but the first letter of each string into the code, then truncates the result to be at most four characters long. Zeros are added at the end if necessary to produce a four-character code. For example, Washington is coded W-252 (W, 2 for the S, 5 for the N, 2 for the G, remaining letters disregarded), and Lee is coded L-000 (L, 000 added). A limitation of Soundex is that it does not know the intricacies of complex spelling rules for English, and because it works on a letter-by-letter

basis, it often does not produce the expected result. Another limitation is that truncating the words to four-character code ignores differences in long strings, which may not be appropriate when finding alternatives for misspelled words. An advantage of Soundex is the small table size and simplicity of the letter-by-letter algorithm, which can provide significant speedup over the other phonetic methods.

Table 1: Soundex coding rules

<i>Code</i>	<i>Letters</i>
0 (not coded)	A, E, I, O, U, H, W, Y
1	B, F, P, V
2	C, G, J, K, Q, S, X, Z
3	D, T
4	L
5	M, N
6	R

1. Replace all but the first letter of *s* by its phonetic code.
2. Eliminate any consecutive repetition of codes.
3. Eliminate all occurrences of code 0 (i.e., eliminate vowels, and the letters *H*, *W* and *Y*).
4. Return the first four characters of the resulting string.

Figure 1: The Soundex algorithm

PHONIX is similar to Soundex in that letters are mapped to a set of codes. Prior to this mapping however, PHONIX applies preliminary transformations to letter groups in order to reduce strings to a canonical form. For example, gn, ghn, and gne are mapped to n, the sequence tjV (where V is any vowel) is mapped to chV if it occurs at the start of a string, and x is transformed to ecs. PHONIX applies altogether about 160 of these transformations. These transformations provide a certain degree of context for the phonetic coding and allow, for example, c and s to be distinguished, which is not possible under Soundex. The Phonix codes are shown in Table 2.

Table 2: PHONIX coding rules

<i>Code</i>	<i>Letters</i>
0 (not coded)	A, E, H, I, O, U, W, Y
1	B, P
2	C, G, J, K, Q
3	D, T
4	L
5	M, N
6	R
7	F, V
8	S, X

The Metaphone algorithm is also a system for transforming words into codes based on phonetic properties. However, unlike Soundex, which operates on a letter-by-letter scheme, Metaphone analyzes both single consonants and groups of letters called diphthongs, according to a set of rules for grouping consonants, and then mapping groups to Metaphone codes.

A drawback of these algorithm as pointed out earlier is that these are language-specific, and typically designed for the English language. There have been attempts to modify these algorithms for other European languages, such as Spanish, Polish and Portugese, but not for Bangla to the best of our knowledge. Using recent research on machine learning methods for letter-to-phoneme conversion [14, 15], application of these techniques to Bangla should be straightforward, provided that there is sufficient training data. As the first step in defining a comprehensive phonetic matching technique for Bangla, we describe an algorithm based on the widely-used Soundex algorithm, suitably modified to reflect Bangla phonetics.

3. Bangla phonetically similar error correction

Bangla letters are partitioned according to phonetic similarity (e.g., I:II, U:UU, NA:NNA, SA:SSA:SHA, etc), with each set represented by a single code. This coding can then be applied to Bangla dictionary to convert it to a non-homophonous one, with each entry pointing to the set of words that correspond to this code [5]. When checking the spelling of a word, we first search for its encoded version in the modified dictionary. If the entry exists, then either the original exists in the list of words corresponding to this code (in which case, the spelling

is correct), or the word is misspelled and the list is offered as the set of alternatives for the original word. If the entry does not exist, then the alternatives must be suggested using one of the edit-distance algorithms (e.g., Levenshtein [15]). However, this technique does not work if an extra error occurs in the spelling, so this technique must be used with a edit-distance algorithm to be effective in a spelling checker. See [2] for a summary of the various commonly-used edit-distance algorithms. This technique works for Bangla conjuncts as well, but only if we eliminate the hasant character from our en-coded strings.

3.1. Soundex algorithm for Bangla

The Soundex algorithm, unmodified, presents a set of difficulties when used in a Bangla spelling checker. In this section, we present some of the prominent issues.

Case 1: Soundex does not consider the first letter in the string.

Problem: This is in fact a general problem with Soundex. If there is a spelling error in the first character of the word, the correct suggestion cannot be produced using Soundex. For example, if we write গুম instead of ঘুম, Soundex will not be able to suggest the correct alternative, as the incorrectly spelled word গুম will begin with গ independent of the character encoding used, Unicode or otherwise at the beginning. Since the phonetically encoded lexicon will have the word ঘুম encoded as something that begins with ঘ, the phonetic method will never produce ঘুম as a suggestion for গুম. Of course, other edit-distance algorithms (e.g., Levenshtein [15]) are able to produce the correct suggestion in this particular case, so a spelling checker employing other similarity measures will produce the expected result (See [2] for a summary of the various edit-distance algorithms).

Case 2: Soundex excludes vowels when encoding strings.

Problem: The অ vowel is often used as a prefix to negate the meaning of Bangla words, and excluding it will often produce suggestions that are of the opposite meaning than the intended one. This may be appropriate behavior for some applications, but not for a spelling checker. For example, the words সূখ and অসূখ will result in the same Soundex code, even though we do not expect one as the suggested alternative for the other, much like we would not expect unwell as the suggested alternative for well .

Problem: Another problem of excluding the vowels is that words that are not phonetically similar and have a very different meanings also produce the same code. বন and বানি, অকাজ and কাজি. বন (forest) and বানি for example will produce the same code if we exclude vowels, where these words do not have same meaning, and in addition, these are phonetically quite different. Similarly, in the case of অকাজ and কাজি , the অ from অকাজ and the ি from কাজি will be excluded to produce the same code, another undesired result.

Case 3: In soundex, consecutive repetitions of the same coded characters are eliminated.

Problem: Unicode specifies that the consonants that make up Bangla juktakkhors are separated by hasant character, which is not coded in our algorithm (i.e., eliminated during the phonetic encoding process). The side-effect of this decision to eliminate hasant is that, at least for a set of juktakkhors, consecutive repetitions of the same consonants will have the same code as the single instance of that con-sonant. Using our algorithm, গ (গ) for example will have the same code as গ, since we exclude the hasant embedded in the Unicode representation of the conjunct. This particular problem is not a general Soundex problem, but rather a consequence of the way our algorithm handles Bangla conjuncts.

3.2. Phonetic matching technique for Bangla

Table 3 shows the proposed Bangla phonetic codes using rules found in [17, 18], and Figure 2 shows the modified Soundex algorithm suitable for a Bangla spelling checker.

Table 3: Bangla phonetic coding rules

Code	Character	Unicode
0 (not coded)	্ (hasant)	09CD
a	আ	0985
	অ	0986
	য়	09DF
	া	09BE
i	ঈ	0987
	ই	0988
	ি	09BF
	ী	09C0

Code	Character	Unicode
u	উ	0989
	ঊ	098A
	ু	09C1
	ূ	09C2
y	এ	098F
	ঈ	0990
	ি	09C7
	ৈ	09C8
o	ও	0993
	ঔ	0994
	ো	09CB
	ৌ	09CC
k	ক	0995
	খ	0996
g	গ	0997
	ঘ	0998
c	চ	099A
	ছ	099B
j	জ	099C
	ঝ	099D
	য	09AF
t	ট	099F
	ঠ	09A0
d	ড	09A1
	ঢ	09A2
n	ন	09A8
	ণ	09A3
f	ত	09A4
	থ	09A5
q	দ	09A6
	ধ	09A7
p	প	09AA
	ফ	09AB

Code	Character	Unicode
b	ব	09AC
	ভ	09AD
m	ম	09AE
	ঙ	0999
	ং	0982
r	র	09B0
	ড়	09DC
	ঢ়	09DD
	ঝ	098B
s	স	09B8
	ষ	09B7
	শ	09B6
h	হ	09B9
	ঃ	0983
l	ল	09B2

1. Replace all of *s* by its phonetic code.
2. Eliminate all occurrences of code 0 (i.e., eliminate *hasant*).
3. Return the resulting string.

Figure 2: The Soundex algorithm for Bangla

3.3. Example of error correction using phonetic matching

Table 4 shows a set of misspelled words, their corresponding encoded versions, and the suggested alternatives.

Table 4: Suggestions for misspelled words

Input	Encoded	Suggestion
খুমাড়	kumar	কুমার
পাসাল	pasan	পাশাণ
দগধ	qgq	দগ্ধ (দগ্ধ)

The case of দগ্ধ deserves special mention. The Unicode representation of the juktakkhor ঙ্গ consists of the letters গ and ধ, separated by a hasant. Since the pronunciation of this sequence of letters is the same

with or without the hasant, it can be safely excluded from the encoding. One the other hand, if we include hasant in the encoding, there will be an one-character error, and the encoded versions will not match, and the correct suggestion can not be produced. A consequence of this handling of the conjuncts is that the large number of Bangla conjuncts does not affect the table, an important consideration when dealing with large dictionaries.

4. Conclusion

We present a preliminary effort at creating a phonetic matching algorithm for Bangla based on Soundex, and tailored for a spelling checker. We describe a prototypical phonetic coding rules and the associated algorithm that produces “good” suggestions for misspelled Bangla words. There are however many complex spelling rules that are not yet addressed in this encoding, such as those involving the use of Reph and Yaphalaa [17, 18], due to the lack of con-text information in our algorithm. The approaches used by PHONIX and Metaphone variants do provide some context, and our future work in this area will concentrate on creating transformation maps to reduce strings to canonical forms before the table-driven encoding step.

5. References

- [1] K. Kukich, “Techniques for automatically correcting words in text”, *Computing Surveys*, 24, (4), 1992, pp. 377–440.
- [2] J. Zobel and P. Dart, “Finding Approximate Matches in Large Lexicons”, *Software - Practice and Experience*, 25(3), March, 1995, pp. 331-345.
- [3] F. Damerau, “A technique for computer detection and correction of spelling errors”, *Communication of the ACM*, 7(3), 1964, pp. 171-176.
- [4] V. Hodge and J. Austin, “A Novel Binary Spell Checker”, *Proc. International Conference on Artificial Neural Networks*, Vienna, August, 2001.
- [5] B. B. Chaudhuri, “Reversed word dictionary and phonetically similar word grouping based spell-checker to Bangla text”, *Proc. LESAL Workshop*, Mumbai, 2001.
- [6] A.B.A Abdullah and A. Rahman, “A Different Approach in Spell Checking for South Asian Languages”, *Proc. 2nd International Conference on Information Technology for Applications (ICITA)*, China, 2004.
- [7] P. Kundu and B.B. Chaudhuri, “Error Pattern in Bangla Text”, *International Journal of Dravidian Linguistics*, 28(2), 1999.
- [8] The Unicode Consortium, *The Unicode Standard, Version 4.0*, Addison-Wesley, 2003. Also available online at <http://www.unicode.org/versions/Unicode4.0.1>.
- [9] D. E. Knuth, *The Art of Computer Programming, Vol. 3*, Addison-Wesley Publishing Company, Reading, Massachusetts, 2nd edition, 1982.
- [10] The Soundex Algorithm, available online at http://www.archives.gov/research_room/genealogy/census/soundex.html.
- [11] L. Phillips, “Hanging on the Metaphone”, *Computer Language*, 7(12), 1990.
- [12] L. Phillips, “The Double Metaphone Search Algorithm”, *C/C++ Users Journal*, 18(6), June, 2000. Also available online at <http://www.cuj.com/documents/s=8038/cuj0006philips/>.
- [13] T. N. Gadd, “PHONIX: The Algorithm”, *Program*, 24(4), pp. 363-366, 1990.
- [14] W. M. Fisher, “A statistical text-to-phone function using n-grams and rules”, *Proc. ICASSP-99, the 1999 IEEE International Conference on Acoustics, Speech and Signal Processing, volume 2*, March 1999, pp 649-652.
- [15] K. Toutanova and R. C. Moore, “Pronunciation modeling for improved spelling correction”, July 2002.
- [16] V. L. Levenshtein, “Binary codes capable of correcting deletions, insertions, and reversals”, *Soviet Physics Doklady*, 10, 1966.
- [17] *Bangla Uccharon Obidhan*, Bangla Academy, Dhaka, Bangladesh.
- [18] *Bangla Banan Obidhan*, Bangla Academy, Dhaka Bangladesh.