

Figure II. System Design

A. Pre-Processing and tokenization

First a text document is read line by line from corpus and each line is pre-processed by elimination of non-Telugu characters, numerals and special characters like colons, semicolons and quotes. Then a pre-processed document is tokenized and extracts the raw words. Words in Telugu text are separated by spaces and are extracted with spaces as delimiter from the document and place all raw words in *Input File*.

Example 1:

Telugu document: ఒక వనంలో గజేంద్రం అనే మదపుటీనుగు తన భార్యలైన దశలక్షకోటి ఆడ ఏనుగులతో విహారిస్తూ దాహం తీర్చుకోడానికి మహావనం మధ్యనున్న సరోవరానికి బయలుదేరింది. దాహం తీరాక జలవిహారంపై బుద్ధిపుట్టి, తన ఆడ ఏనుగులతో సరోవరాన్ని అల్లకల్లోలం చేస్తూ విజృంభించి జలకేళిలో లీనమై ఉన్న గజేంద్రాన్ని, గొప్ప మొసలి కోరలు గుచ్చి ముందరి కాలొకటి పట్టుకొన్నది.

Pre-processed Telugu document is

ఒక వనంలో గజేంద్రం అనే మదపుటీనుగు తన భార్యలైన దశలక్షకోటి ఆడ ఏనుగులతో విహారిస్తూ దాహం తీర్చుకోడానికి మహావనం మధ్యనున్న సరోవరానికి బయలుదేరింది దాహం తీరాక జలవిహారంపై బుద్ధిపుట్టి తన ఆడ ఏనుగులతో సరోవరాన్ని అల్లకల్లోలం చేస్తూ విజృంభించి జలకేళిలో లీనమై ఉన్న గజేంద్రాన్ని గొప్ప మొసలి కోరలు గుచ్చి ముందరి కాలొకటి పట్టుకొన్నది

Tokenized Document is

ఒక వనంలో గజేంద్రం అనే మదపుటీనుగు తన భార్యలైన దశలక్షకోటి ఆడ ఏనుగులతో విహారిస్తూ దాహం తీర్చుకోడానికి మహావనం మధ్యనున్న సరోవరానికి బయలుదేరింది దాహం తీరాక జలవిహారం పై బుద్ధిపుట్టి తన ఆడ ఏనుగులతో సరోవరాన్ని అల్లకల్లోలం చేస్తూ విజృంభించి జలకేళిలో లీనమై ఉన్న గజేంద్రాన్ని గొప్ప మొసలి కోరలు గుచ్చి ముందరి కాలొకటి పట్టుకొన్నది

All these tokens are stored in *Input File*.

B. Pseudo N-gram and categorization

Pseudo N-gram is a procedure that reduces words by stripping derivational and inflectional suffixes from each word to get valid root[1]. It takes raw words from *Input File* as input. Read one word at a time from file. First find the length of word and fix the stripping length. Stripping length will be varied based on word length. Maximum stripping length is 5 and minimum is 2, and then apply Pseudo N-gram algorithm for each word. For each step, strip the word from end and check, if it is valid root or not. If it is valid root, then extract root, accept and perform categorization is shown in fig.3 If it is not a valid root, decrease the stripping length by one and check. This process is repeated until stripping length not equals to zero.

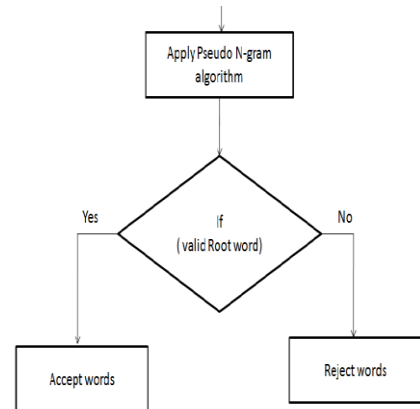


Figure III. Pseudo N-gram

Document categorization using words is complex task due to the nature of the corpus. To reduce the computational complexity it is necessary to adopt a corpus database with base words. Pseudo N-gram based categorization of words is one of the solution to achieve the above goal. The proposed algorithm is as follows.

Pseudo N-gram Algorithm

1. Start
2. Take the LIST [1000] [12] of words as input.
3. SET i=0
4. WHILE (LIST[i] != NULL) Repeat the steps from 5 to 18
 Otherwise go to step 19
5. Read one WORD at a time from LIST i.e WORD = LIST[i]
6. Find the length of WORD as Word_Len=strlen(WORD)
7. SET STRIPPING_LENGTH=0
8. IF (Word_Len > 2) THEN go to step 9 otherwise go to step 18
9. IF (Word_Len >=7)
 THEN
 STRIPPING_LENGTH=5
 ELSE
 IF (Word_Len < 7 OR Word_Len >=5)
 THEN
 STRIPPING_LENGTH=4
 ELSE
 STRIPPING_LENGTH=3
 END IF
10. SET Count=0
11. WHILE (STRIPPING_LENGTH=>0) Repeat steps from 12 to 17
12. IF(Count<
 Word_Len- STRIPPING_LENGTH- 1)
 THEN repeat steps 13 and 14

13. TEMP_WORD [Count] = WORD[Count]
14. SET Count= Count+1 then go to Step 12
15. IF (TEMP_WORD == Valid Root WORD)
 /*Validate with Telugu dictionary */
 THEN Go to Step 16 otherwise Go to step 17
16. WRITE WORD in Accepted File Go to step 18
17. SET STRIPPING_LENGTH= STRIPPING_LENGTH – 1 Go to Step 11
18. SET i = i+1 then go to step 4
19. EXIT

III. TESTING AND RESULTS

The experiments were conducted on Telugu Corpus, collected from online newspapers and wikipedia. This work has been implemented on sample selection of 1,550 documents. A sequence of words from the word lists was used in extracting valid root by pseudo N-gram algorithm and the results are presented in Table 1, which contains list of words with initial stripping and final root word for evaluation.

Example 2: List[20][10]={ పాలసముద్రంలో, ఏనుగులతో, మదపుటీనుగు, భార్యలైన, జలవిహారంపై, సరోవరానికి, నిలబడ్డాయిగాని, తోచలేదు, తప్పించాలి, అల్లకల్లోలం, బుద్ధిపూట్టి, చెల్లాచెదురుగా }

TABLE 1.
Results of Pseudo N-gram algorithm

List of Words before pseudo N-gram	Initial Word length	Initial Stripping Length	Final stripping length to make a Valid Word	Stripped Suffix	Valid Root word
పాలసముద్రంలో	6	4	1	లో	పాలసముద్రం
ఏనుగులతో	5	4	2	తో	ఏనుగు
మదపుటీనుగు	6	4	0	---	మదపుటీనుగు
భార్యలైన	4	3	2	లైన	భార్య
జలవిహారంపై	6	4	1	పై	జలవిహారం
సరోవరానికి	6	4	0	---	Not a Valid Root
నిలబడ్డాయిగాని	7	5	2	గాని	నిలబడ్డాయి
తోచలేదు	4	3	0	---	తోచలేదు
తప్పించాలి	4	3	1	లి	తప్పించా
అల్లకల్లోలం	5	4	0	---	అల్లకల్లోలం
బుద్ధిపూట్టి	4	3	2	ట్టి	బుద్ధి
చెల్లాచెదురుగా	6	4	1	గా	చెల్లాచెదురు

When the words list was given to pseudo N-gram there were some set of words that could not be recognized as valid root. In the above table, the word సరోవరానికి is not a valid root. After applying the pseudo N-gram algorithm, the possible words are సరో, సరోవ, సరోవరా, సరోవరాని, and సరోవరానికి for each stripping వరానికి, రానికి, నికి, and కి with the stripping length like 4,3,2,1 and 0. Valid root for సరోవరానికి is సరోవరం. In proposed algorithm extraction of valid root using pseudo N-gram based categorization are addressed. Disadvantage of pseudo N-gram is, it generates some words that are not valid.

IV. CONCLUSION

Extraction of valid root is already proven method of Information Retrieval in Indian languages. Pseudo N-gram technique is a new language independent and they are well suited for different complex Indian languages like Hindi and Kannada. In this paper, accuracy of pseudo N-gram is more than N-gram model. The maximum accuracy observed is 85% for pseudo N-gram. There is no such report of Extracting Telugu language valid root words using pseudo N-gram. As part of our research work in Telugu categorization, we propose to extend it for recognize all words as valid roots.

Acknowledgements



N. Swapna received B.Tech. in Computer Science and Information Technology from VREC-Nizamabad, JNT University. M.Tech. in Computer Science & Engineering from JNTU Anantapur and she is Pursuing Ph.D. in the area of Information Retrieval Systems from the Department of Computer Science and Engineering, JNT University Hyderabad. She has 12 years of teaching experience in various engineering colleges. Currently she is working as Associate Professor and Training & Placement Officer in the department of Computer Science Engineering ,Vijay Rural Engineering College, Nizamabad, India. To her credit Mrs.Swapna Narala has 10 publications in various National / International Conference and Journals. She is also a Member of Various Technical Bodies including IEEE, ISTE etc. Her area of interest includes Information Retrieval, Text Mining, Web Mining, Machine Learning, Information Security etc.



B. Padmaja Rani received B.Tech Electronics Engineering from Osmania University, M.Tech in Computer Science from JNT University Hyderabad, India and she has been awarded Ph.D. in Computer Science from JNT University, Hyderabad, India. At present she is working as Professor in the Department of Computer Science and Engineering, JNTUH College of Engineering, JNTU University Hyderabad. She is having 20 years of experience in Industry and Academia. At present she is a Professor of Computer Science and Engineering Department in JNTUH College of Engineering, JNT University, Hyderabad. Her area of Research includes Information Retrieval, Data Mining, Machine Translation, Computer Networks, Software Engineering etc. She is guiding 6 Research Scholars in the area of Information Retrieval and Computer Networks. To her credit she is having more than 60 publications in reputed International Journals and Conferences. She is a member of various advisory committees and Technical Bodies. She is also a Member of Various Technical Associations including ISTE, CSI, IEEE etc.

REFERENCES

- [1] Porter M.F. " An algorithm for suffix stripping",pogram,14(3),1980,130-137
- [2] Uppal sharma," unsupervised learning of morphology of a highly inflectional language" Ph.D thesis submitted to Department of Computer Science and Information Technology",Tezpur University, Napaam,Assam,India,2006.
- [3] A D Manning, P. Raghavan, and H. Scutze., An introduction to information retrieval, Cambridge: Cambridge university press, Vol. 1, 2009, PP:6.
- [4] U.Rao, 2008, Functional Specifications of Morphology CLATS, Hyderabad Central University, Version 1.3.1, 2008, PP:1-32.
- [5] Bharadwaja Kumar, G., Kavi Narayana Murthy, and B. B. Chaudhuri. , Statistical analyses of Telugu text corpora, International journal of Dravidian linguistics (IJDL), vol-36, issue- 2, 2007, PP:71-99.
- [6] Mrs.A.Kanaka Durga, Dr.A.Govardhan, Ontology Based Text Categorization Telugu Documents, International Journal of Scientific & Engineering Research Volume 2, Issue 9, September-2011 , PP: 1-4
- [7] Murthy, Kavi Narayana, and G. Bharadwaja Kumar., Language identification from small text samples, Journal of Quantitative Linguistics, vol-13, issue-1, 2006, PP:57-80.
- [8] Rani, Dr B. Padmaja, and Dr A. Vinay Babu. , Novel Implementation of Search Engine for Telugu Documents with Syllable N-Gram Model, International Journal of Engineering Science and Technology, vol- 2, issue-8, 2010, PP:3712-3720.
- [9] Kalyani, N. and Sunitha, K.V.N. (2009): A Novel approach to Improve rule based Telugu Morphological Analyzer, World Congress on Nature & Biologically Inspired Computing (NaBIC 2009).