Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

Telugu Text Categorization using Language Models

K. Ramakrishna

Received: 6 December 2015 Accepted: 5 January 2016 Published: 15 January 2016

5 Abstract

1

2

3

18

6 Document categorization has become an emerging technique in the field of research due to the

7 abundance of documents available in digital form. In this paper we propose language

8 dependent and independent models applicable to categorization of Telugu documents. India is

⁹ a multilingual country; a provision is made for each of the Indian states to choose their own

¹⁰ authorized language for communicating at the state level for legitimate purpose. The

¹¹ availability of constantly increasing amount of textual data of various Indian regional

¹² languages in electronic form has accelerated. Hence, the Classification of text documents

¹³ based on languages is crucial. Telugu is the third most spoken language in India and one of

¹⁴ the fifteen most spoken language n the world. It is the official language of the states of

Telangana and Andhra Pradesh. A variant of k-nearest neighbors algorithm used for
 categorization process. The results obtained by the Comparisons of language dependent and

16 categorization process. The results obtained by the Comparisons of la 17 independent models.

Index terms— text categorization, language dependent and independent models, k-nearest neighbors. 19 I. Introduction ow a day's huge amount of information is being posted on to the web. In order to get useful 20 information from the web, the information available has to be categorized. Text Categorization is the task of 21 automatically categorizing a set of unlabeled text documents to their corresponding categories from a predefined 22 category set [2]. These categories can be viewed as a set of documents and test document can be treated as a 23 query to the system. The measures to evaluate the information retrieval systems are often applicable to measure 24 25 effectiveness text categorization systems [1]. Text categorization has many applications [2], like information retrieval system, search engine, text filtering, word sense disambiguation, language identification, POS tagging 26 and machine translation etc. Telugu is one of the old and traditional languages of India and it is categorized 27 as one of the Dravidian language family unit with its own high-class script. It is the authorized language of 28 the Telangana and Andhra Pradesh states in south India. Amit et al ??6] surveyed that in India the Telugu 29 native speakers are above 50 million. It was positioned between 13 to 17 largest spoken languages all over the 30 world. Telugu is a rich morphological language that has high word conflation ???]. Various approaches for text 31 categorization have been done on Indian languages. Most of the works have been reported on Telugu language. 32 M Narayana Swamy et al have used KNN, NB and decision tree classifier [4]. They have experiment on Kannada, 33 Tamil and Telugu corpus statistics is illustrated by Zipf's law. Analysis of N-gram model on text classification 34 was proposed in the work of [5]. Goverdhan. A Durga k et al [3] projected a technique with ontology text 35 categorization for Telugu digital-items and retrieval system. For the best of our knowledge, this is the first time 36 our proposed language models have been applied for Telugu text categorization. The paper is structured as 37 follows; section 2 describes the system overview, section 3 explains Testing and results and at the last, a section 38 4 conclusion is drawn. 39

40 1 II. System Overview

41 The system design of the proposed approach can be shown in the Figure ??1. First read a text document from 42 corpus and each line is pre-processed by elimination of non-Telugu characters, numerals and special characters 43 like colons, semicolons and quotes. Then a pre-processed document is tokenized and extracts the raw words. 44 Words in Telugu text are separated by spaces and are extracted with spaces as delimiter from the document 45 and place all raw words in Input File. Language dependent and independent models are takes raw words from 46 Input File as input. Read one word at a time from file. Finally find the root word by applying various models 47 like vibhaktulu based stemming, suffix removal stemming, Rule based suffix removal stemming, N-gramming,

a pseudo N-gramming and Rule based Pseudo N-gramming. Finally, apply the text categoryzation. Our proposed

⁴⁹ language models are categorized in three ways are shown in figure 2. These models take raw words from Input

 $\,$ 50 $\,$ File as input and identify the root word.

⁵¹ 2 c) Suffix removal stemming

Suffix removal stemming is the process of finding the root word from the word by removing the matched suffix 52 with suffix list which is shown in figure 3. By observing the Telugu data set, it is found that maximum suffix 53 length will be 2(two) and minimum is one. Suffix removal stemming method giving better performance than 54 vibhaktulu based stemming algorithm. It's accuracy is 58-59%. Suffix removal stemming is a base method for 55 Rule based Suffix removal stemming algorithm. The result of suffix removal stemming words may normally 56 contain inflections. The inflections in the stem word cannot be removed using simple suffix removal. We have 57 designed rule based suffix removal of some possible inflections that frequently occur in the Telugu Language. The 58 rules are used to replace characters are presented in Table ??. By these rules the electiveness of the proposed 59 Rule based Suffix removal stemming algorithm is increased. Accuracy of Rule based suffix removal is 69-70%. 60

⁶¹ 3 Table 1: Rules for Replacement Syllables e) Pseudo N-

62 gramming

Pseudo N-gram is the process of finding the root word by stripping the word from the end. Stripping length will
be taken depending on the word length. Maximum stripping length is 5 and minimum is 2. Example of Pseudo
N-gramming is shown in figure 4. It is a language independent.

66 4 f) Rule Based Pseudo N-Gramming

It is a hybrid model. Pseudo N-gram is a base method for this processing to remove suffixes from words. The
result of Pseudo N-gram of some words normally contains inflections. The inflections in the stem word cannot
be removed using simple Pseudo N-gram.

We have designed rule based Pseudo N-gram which contain set of rules used to replace characters. These rules used for words normally contain more inflections that frequently occur in the Telugu Language. List of rules with sample example are shown in Table ??.

Table 3: List of rules for Rule based pseudo N-gramming g) K-NN Classifier

The k-NN classifier is a similarity-based learning method that has been shown to be very effective for a variety of problem domains including text categorization [9, 10]. Given a test document, the k-NN method finds the k nearest neighbors among the training documents, and uses the categories of the k neighbors to weight the

category. The similarity score of each and every neighbor document to the test document is used as the weightof the classes of the neighbor document.

80 6 III. Testing and Results

The proposed models are evaluated on Telugu Corpus, collected from online newspapers and Wikipedia. This work has been implemented on sample selection of 1,500 documents of seven categories are presented in Table ??.

⁸⁴ 7 Table 4: Categories of Telugu Documents

To evaluating the performance of the proposed system using KNN classification, we use the typical evaluation metrics that come from information retrievalprecision (P), recall (R), and F1 measure: Where TP is True Positives, TN is True Negatives, FN is False Negatives and FP is False Positive [8]. We have projected the

performance of the proposed language models result with KNN classifier shown in Table ??.

89 8 IV. Conclusion

⁹⁰ In this paper, we proposed various language dependent and independent models. Among these models the ⁹¹ performance of Rule based pseudo Ngramming is more. So it is well suited for Telugu Text categorization. As

part of our research work in Telugu categorization, it is also suitable for other complex Indian languages like

⁹³ Hindi, Malayalam and Kannada.

 $^{^{1}}$ © 2016 Global Journals Inc. (US)



Figure 1: Figure 1 :







Figure 3: Figure 3 :

S.No	List of characters/syllable sound found as suffix	Replacement characters ଡ୦, ଭ, ବ	
1	ಅ,ಆ		
2	യ,യറ	డ, అ, అం	
3	a, a + w, & +w	ఇ, అం	
4	ఎ ,ఏ, ఎం	କ୍ ,ଣ ,ଖ, ଡ୦	
5	ఓ ఓ ఔ	a, ¢	
6	అం	e	

Figure 4: Figure 4 :



Figure 5:

List of Words	Intial	Intial	Final	Stripped	Valid
before	Word	Stripping	stripping	Suffix	Root word
pseudo N-gram	length	Length	length to		
			make a		
			Valid Word		
పాలసముద్రంలో	6	4	1	లో	పాలసముద్రం
ఏనుగులతో	5	4	2	లతో	ఏనుగు
మదపుటేనుగు	6	4	0		మదపుటేనుగు
భార్యలైన	4	3	2	లైన	భార్య
జలవిహారంపై	6	4	1	ي ک	జలవిహారం
సరోవరానికి	6	4	0		Not a Valid
					Root
నిలబడ్డాయిగాని	7	5	2	గాని	నిలబడ్డాయి
తోచలేదు	4	3	0		తోచలేదు
తప్పించాలో	4	3	1	లో	తప్పించా
అల్లకల్లోలం	5	4	0		అల్లకల్లోలం
బుద్ధిపుట్టి	4	3	2	పుట్టి	బుద్ధి
చెల్లాచెదరుగా	6	4	1	ന	చెల్లాచెదరు

515

Figure 6: Table 5 : 1 HFigure 5

S.No	List of characters/syllable sound found as suffix	Replacement characters	List of Words are not recognized by Pseudo N- gram	List of words recognized by Rule based Pseudo N-gram
1	ಅ,ಆ	అం, టి, ఇ	త ప్పడానికి ప్రమాదాన్ని కరటాల వర్షానికి గంపడాశ పళ్ళయిన	చెప్పడం ప్రమాదం కెరటం వర్షం గంపడు పళ్ళి
2	డి 'దం	ఉ, అ, అం	నిపిదించడం దేవుడి హింసించే	నిషిదం దేవుడు హింస
3	డ , డ + లు , ఓ +లు	യ, ಅಂ	ఆక్యరుల్ని ఎడారులు సన్నజాజాల చీకట్	ఆక్యతి ఎదారి సన్నబాబి చీకటి
4	ఎ ,ఏ, ఎం	ఇ ,డ ,అ, అం	చోటిక్కడ మురిపెంగ పసేమిలేదు ఎక్కడెక్కడో	చోటు మురిపం పని ఎక్కడ
5	ఓ ఓ ఔ	ద, కా	కాలొకటి రాడీపెడో	కాలు రాడు
6	ക്ര	â	పగలంతా క ళ్ళం తా	పగలు కళ్ళు

		1			
110	120	247	100	268	80
10,640	88,427	26,255	87,552	99,964	27,061
	110 10,640	110 120 10,640 88,427	110 120 247 10,640 88,427 26,255	110 120 247 100 10,640 88,427 26,255 87,552	110 120 247 100 268 10,640 88,427 26,255 87,552 99,964

Figure 8:

 $\mathbf{2}$

Figure 9: Table 2 :

- [Landauer et al. ()] 'An Introduction to Latent Semantic Analysis'. T K Landauer , P W Foltz , D Laham .
 Discourse Processes, 1998. p. .
- ⁹⁶ [Vishnu Vardhan ()] Analysis of N-gram model on Telugu Document classification thesis, B Vishnu Vardhan .
 ⁹⁷ 2008.
- ⁹⁸ [Murthy ()] 'Automatic Categorization of Telugu News Articles, Department of Computer and Information
 ⁹⁹ Sciences'. K N Murthy . *Hyderabad*, 2003. University of Hyderabad
- 100 [Indian Language Text Representation and Categorization Using Supervised Learning Algorithm M Narayana Swamy1]
- Indian Language Text Representation and Categorization Using Supervised Learning Algorithm M Narayana
 Swamy1,
- 103 [Mrs et al. (2011)] 'Ontology Based Text Categorization Telugu Documents'. A Mrs, Kanaka Durga, . A Dr,
- 104 Govardhan . International Journal of Scientific & Engineering Research September-2011. 2 (9) p. .