

# SENTIMENT ANALYSIS OF NEPALI SENTENCES USING NAIVE BAYES CLASSIFICATION

ASHOK CHHETRI (7926/072)

ABHISHEK SAPKOTA (7921/072)

MAHESH ACHARYA (7942/072)

RABIN BHANDARI (7950/072)

ASIAN SCHOOL OF MANAGEMENT AND TECHNOLOGY

UNDER THE SUPERVISION OF

MR. BIKASH BALAMI

# TABLE OF CONTENT

- . INTRODUCTION
- . PROBLEM STATEMENT
- . OBJECTIVES
- . SCOPE AND LIMITATIONS
- . METHODOLOGY
- . PROJECT SCHEDULE
- . EXPECTED OUTCOME
- . REFERENCES



# INTRODUCTION

- . SENTIMENT ANALYSIS IS CONTEXTUAL MINING OF TEXT
- . IT SPECIFIES THE POLARITY (POSITIVE, NEGATIVE, AND NEUTRAL) OF THE SENTENCES
- . SENTIMENT ANALYSIS IS DIFFICULT AND CHALLENGING FOR INDO-ARYAN LANGUAGE IN WHICH MACHINE LEARNING APPROACHES ARE USED IN SENTIMENT ANALYSIS BY ANALYZING HUGE AMOUNT OF DATA
- . THERE ARE DIFFERENT ALGORITHM ARE USED IN CLASSIFICATION OF SENTIMENT LIKE SVM AND NAÏVE BAYES
- . AMONG THEM WE USE NAÏVE BAYS CLASSIFICATION ALGORITHM

# PROBLEM STATEMENT

SENTIMENT ANALYSIS IS IN RESEARCH PHASE IN ARYAN INDO LANGUAGE; NEPALI LANGUAGE. DUE TO LACK OF RESOURCES SENTIMENT ANALYSIS IS DIFFICULT AND CHALLENGING FOR OUR LANGUAGE

IN THIS PROJECT WORK , PROBLEM OF NEPALI SENTIMENT IS ADDRESSED

THE RECOGNIZATION TASK IS CARRIED OUT BY SUPERVISED MACHINE LEARNING NAÏVE BAYES

.EG

मलाई यो मन परेन।[ NEGATIVE ]

आहा कती राम्रो घडी।[ POSITIVE ]



# OBJECTIVES

THE MAIN OBJECTIVE OF THE SENTIMENT ANALYSIS SYSTEM IS:

- . TO CLASSIFY THE OPINIONS OF THE CUSTOMERS FOR E-COMMERCE BASED PORTAL.

## LITERATURE REVIEW

THERE ARE LOTS OF RESEARCHES WHICH HAVE BEEN DONE IN THE FIELD OF SENTIMENT ANALYSIS,

IN [3] "SEMANTIC ANALYSIS ON NEPALI MOVIE REVIEWS" AUTHOR ASHOK PANTA NAIVE BAYES SENTIMENT ANALYSIS. THIS SYSTEM USES NAIVE BAYES ALGORITHM UNDER THE SUPERVISED MACHINE LEARNING TO CLASSIFY THE NEPALI MOVIE REVIEWS.

IN [1] "NAMED ENTITY RECOGNITION FOR NEPALI TEXT USING SUPPORT VECTOR MACHINE," AUTHOR SURYA BAHADUR BAM. THIS SYSTEM USED SUPERVISED MACHINE LEARNING TO CLASSIFY NEPALI PHRASES.

IN [2] "SEMANTIC ORIENTATION APPLIED TO UNSUPERVISED CLASSIFICATION OF REVIEWS" AUTHOR PETER D. TURNEY. THIS SYSTEM USED UNSUPERVISED MACHINE LEARNING TECHNIQUE FOR THE CLASSIFICATION OF THE REVIEWS.



# SCOPE AND LIMITATIONS

## SCOPES

- . IT IS USED TO FIND THE OPINION OF THE CUSTOMERS ON THE NEWLY LUNCHED PRODUCTS.
- . IT IS USED IN THE FILM INDUSTRY TO ANALYSIS THE COMMENTS AND REVIEWS OF THE MOVIES AND SONGS.
- . IT IS USED IN NEWS PORTALS TO FIND OUT THE PUBLIC OPINIONS AND REVIEWS.

## LIMITATIONS

- . THE ACCURACY OF THE REVIEWS IS LOW WHEN THERE IS THE MIXTURE OF NEPALI AS WELL AS ENGLISH LANGUAGE IN COMMENTS.
- . AMBIGUITY CAN OCCUR, AND SEMANTIC ANALYSIS IS DIFFICULT TO ADDRESS.

# METHODOLOGY

## DATA COLLECTION

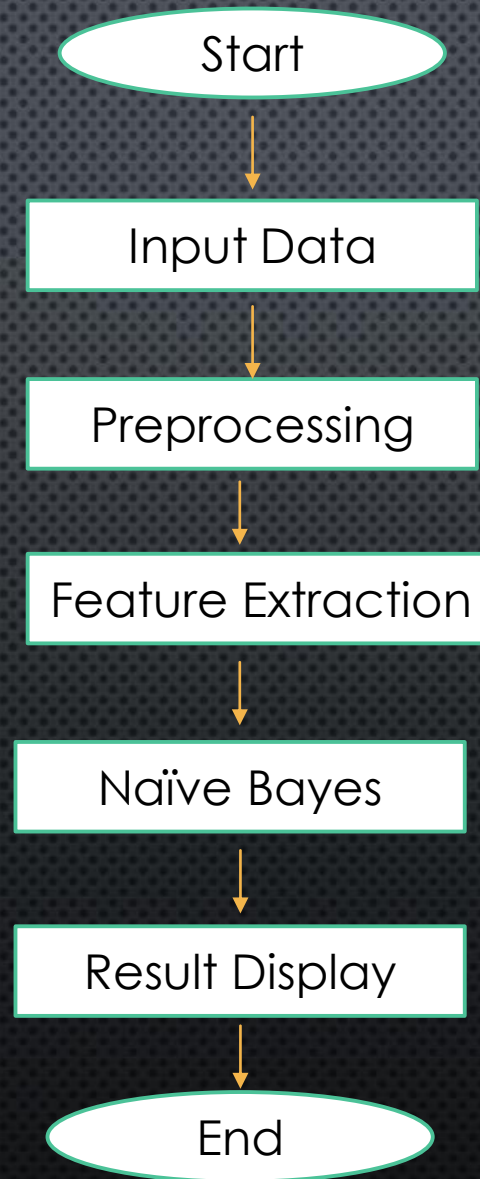
WE ARE MANUALLY COLLECT DATA FROM DIFFERENT NEW AND ECOMMERCE SITES OF NEPAL.

## PREPROCESSING

PREPROCESSING IS USUALLY DONE MANUALLY BY THE HUMANS IN REGARD TO THE DATA SETS. HERE SUPERVISED MACHINE



# PROPOSED IMPLEMENTATION MODEL FOR SENTIMENT ANALYSIS



# FEATURE EXTRACTION

FEATURES IN THE CONTEXT OF OPINION MINING ARE THE WORDS, TERMS OR PHRASES THAT STRONGLY EXPRESS THE OPINION AS POSITIVE OR NEGATIVE

## TF-IDF

FEATURE REPRESENTS WEIGHT OF THE PARTICULAR TERM PRESENT IN THE TEXT DOCUMENT.

$$W_{ik} = \frac{tf_{ik} \log\left(\frac{N}{n_k}\right)}{\sum_{k=1}^t (tf_{ik})^2 [\log\left(\frac{N}{n_k}\right)]^2}$$

WHERE,

- TF = TERM FREQUENCY,
- IDF = INVERSE DOCUMENT FREQUENCY,
- $T_k$  = TERM K IN DOCUMENT  $D_i$ ,
- N = TOTAL NUMBER OF DOCUMENT IN COLLECTION OF C,
- $n_k$  = THE NUMBER OF DOCUMENT IN C THAT CONTAIN T-K,
- $idf_k$  = INVERSE DOCUMENT FREQUENCY OF  $T_k$  IN DOCUMENT ,
- $idf_k = \log\left(\frac{N}{n_k}\right)$



## PRESENCE OF POLAR WORD

POLAR WORDS ARE THE WORDS WHICH REPRESENT THE SENTIMENT LIKE GOOD AND BAD

. EXAMPLE: “यो सामान राम्रो छ।”

## COUNT OF POSITIVE WORD

WE CALCULATED THE NUMBER OF POSITIVE WORDS IN THE SENTENCE AND ADDED IT AS A FEATURE

FOR EXAMPLE, “यो सामान राम्रो छ तेसैले मेरो छोरा ले धेरै मन परायो।”

HAS FOUR POSITIVE WORDS SO IT IS A POSITIVE SENTENCE.

## COUNT OF NEGATIVE WORD

WE ALSO CALCULATED THE NUMBER OF NEGATIVE WORDS PRESENT IN THE SENTENCE AND ADDED IT AS A FEATURE.

E.G. “नेपालमा अन्लाईन सामानमा धेरै ठगि हुने हुँदा मान्छेहरु समान किन्न मन पराउनुन्।”

HAS TWO NEGATIVE WORDS.



# NAIVE BAYES

NAIVE BAYESIAN CLASSIFIER IS A SIMPLE PROBABILISTIC CLASSIFIER BASED ON BAYES THEOREM WITH STRONG INDEPENDENCE ASSUMPTIONS OF FEATURE SPACE. DEPENDING ON THE PRECISE NATURE OF THE PROBABILITY MODEL, NAIVE BAYES CLASSIFIER CAN BE TRAINED VERY EFFICIENTLY IN A SUPERVISED LEARNING SETTING[4].

$$P(H/X) = \frac{P(X/H)P(H)}{P(X)}$$

WHERE,

$P(X/H)$  = POSTERIOR PROBABILITY OF H CONDITIONED X

$P(H/X)$  = POSTERIOR PROBABILITY OF X CONDITIONED H

$P(H)$  = PRIOR PROBABILITY OF HYPOTHESIS H

$P(X)$  = PRIOR PROBABILITY OF X

# CLASSIFY NAIVE BAYES TEXT

POSITION = ALL WORD IN DOC THAT CONTAIN TOKENS FOUND IN VOCABULARY

RETURN VNB ,

WHERE,

$$VNB = \underset{V_J \in V}{\operatorname{ARGMAX}} P(V_J) \prod_{I \in \text{POSITIONS}} P(A_I / V_J)$$

.EG

TEST = मेरो नाम राम्रो छ ।

$$P(\text{TEST}/+) = P(+ )P(\text{मेरो}/+)P(\text{नाम}/+)P(\text{राम्रो}/+)P(\text{छ}/+)$$



## PRECISION

PRECISION (ALSO CALLED POSITIVE PREDICTIVE VALUE) IS THE NUMBER OF CORRECTLY CLASSIFIED POSITIVE EXAMPLES DIVIDED BY THE NUMBER OF EXAMPLES LABELED BY THE SYSTEM AS POSITIVE.

## RECALL

RECALL (ALSO CALLED SENSITIVITY) IS THE NUMBER OF CORRECTLY CLASSIFIED POSITIVE EXAMPLES DIVIDED BY THE NUMBER OF POSITIVE EXAMPLES IN THE TEST DATASET.

## **F-MEASURE**

HARMONIC MEAN OF PRECISION AND RECALL. MATHEMATICALLY,

$$F = \frac{2(P*R)}{P+R}$$

# Project Schedule

[illegible]



## EXPECTED OUTCOME

AT THE END OF THIS PROJECT, THE PROPOSED MODEL IS TO BE EXPECTED TO CLASSIFY GENERAL CUSTOMER OPINION TOWARDS THE PRODUCT IN NEPALI TEXT.

# REFERENCES

- [1]SURYA BAHADUR BAM, "NAMED ENTITY RECOGNITION FOR NEPALI TEXT USING SUPPORT VECTOMACHINE," TRIBHUVAN UNIVERSITY DEPARTMENT OF SCIENCE AND TECHNOLOGY, 2011.
- [2]PETER D. TURNEY, "THUMBS UP OR THUMBS DOWN? SEMANTIC ORIENTATION APPLIED TO UNSUPERVISED CLASSIFICATION OF REVIEWS," INSTITUTE FOR INFORMATION TECHNOLOGY NATIONAL RESEARCH COUNCIL OF CANADA OTTAWA, ONTARIO, CANADA, K1A 0R6, 2002.
- [3]ASHOK PANTA, "SENTIMENT ANALYSIS ON NEPALI MOVIE REVIEWS USING MACHINE LEARNING," TRIBHUWAN UNIVERSITY DEPARTMENT OF SCIENCE AND TECHNOLOGY, 2013.
- [4]LINA L. DHANDE AND DR. PROF. GIRISH K. PATNAIK, —ANALYZING SENTIMENT OF MOVIE REVIEW DATA USING NAIVE BAYES NEURAL CLASSIFIERII, IJETTCS, VOLUME 3, ISSUE 4 JULY-AUGUST 2014, ISSN 2278-6856.
- [5]RAJARAMAN, A., & ULLMAN, J. (2011). DATA MINING. IN MINING OF MASSIVE DATASETS (PP. 1-17). CAMBRIDGE: CAMBRIDGE UNIVERSITY PRESS.



**THANK  
YOU!**

