# University of Science and Technology of Hanoi

## Text Classification

## Using Decision Tree and Maximum Entropy

Student: **Bui Dinh Duong**

Hanoi, September 20th, 2013

# Organization

➢ Introduction

➢ Objectives

➢ Text Classification Overview

➢ Decision Tree

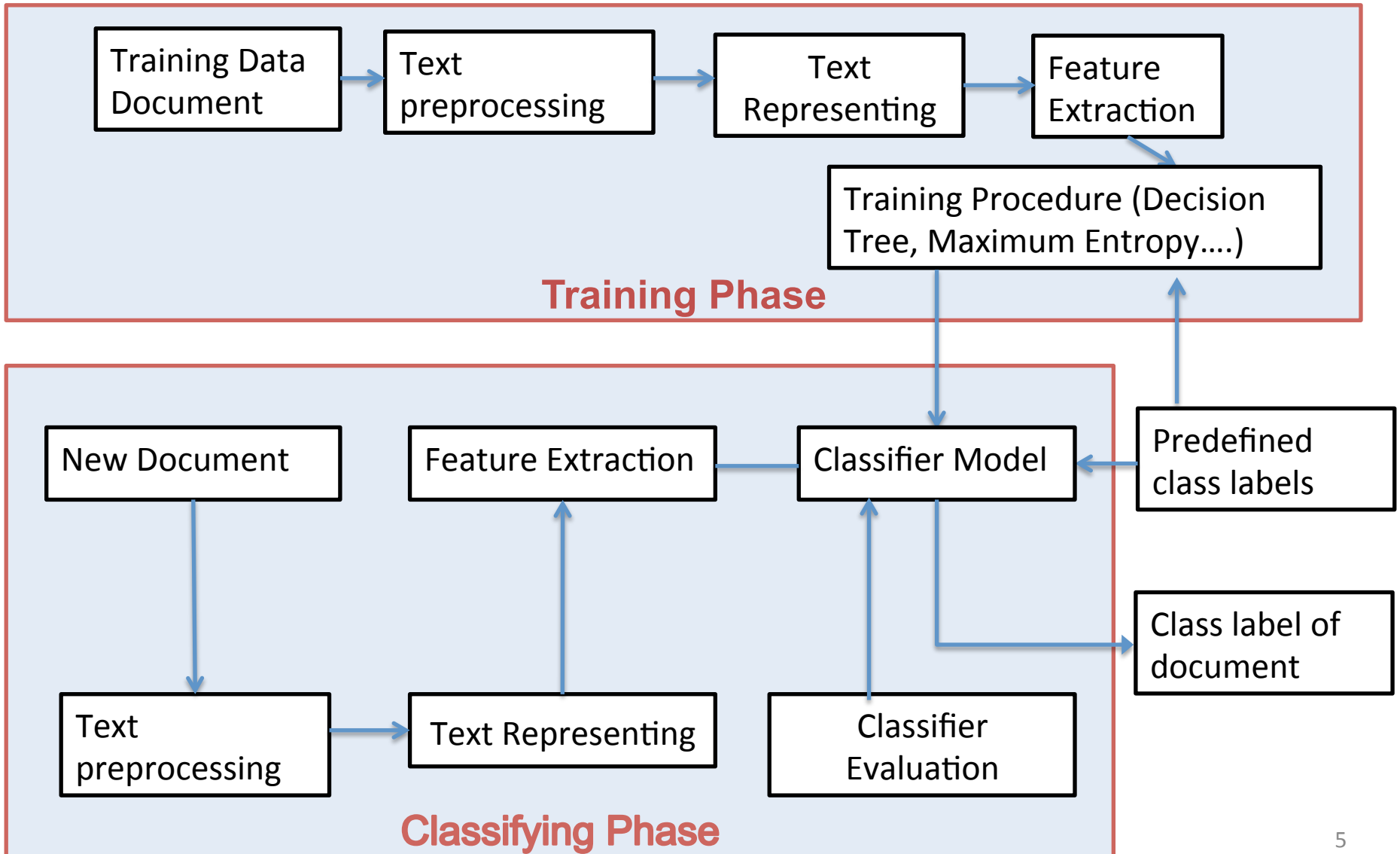➢ Maximum Entropy

➢ Experiment

➢ Conclusion

# Introduction

- ➢ Text classification
  - ▪ Assign a document to one or more predefined classes
- ➢ Applications
  - ▪ E-mail spam filtering
  - ▪ Categorize newspaper articles into topics
  - ▪ Organize Web pages into hierarchical categories
  - ▪ Language identification
- ➢ Methods
  - ▪ Naive Bayes
  - ▪ Maximum Entropy
  - ▪ Decision Tree
  - ▪ Support Vector Machine (SVM)

# Objectives

➢ Study the stages of text classification

➢ Study Decision Tree method

➢ Study Maximum Entropy method

➢ Do experiment in text classification using Weka

# Text Classification

## Training Phase

| Training Data Document | → | Text preprocessing | → | Text Representing | → | Feature Extraction |
|---|---|---|---|---|---|---|

Training Procedure (Decision Tree, Maximum Entropy….)

**Training Phase**

## Classifying Phase

| New Document | Feature Extraction | Classifier Model | Predefined class labels |
|---|---|---|---|

| Text preprocessing | Text Representing | Classifier Evaluation | Class label of document |
|---|---|---|---|

**Classifying Phase**

5

# Text Preprocessing

➤ What is the objective?

- Reduce the size of data
- Get only things we need

➤ How to do?

- Convert document to lower case
- Remove words that rely occur in the document
- Remove special character
- Remove stop-words (words are not used to classify)
- Remove suffix, prefix of word to get the root word ("clusters", "clustering", "clustered" => cluster)

# Text Representing

➢ What is the objective?

- Represent text data in a suitable model to process

➢ How to do?

- Vector Space Model (most popular method)
  - Each document is represented as a vector of **word weighting**

  For example: "**The brown fox jumps over the lazy dog**"

a  an  …brown,..        dog  …      fox      jump      lazi      over the

$( 0, 0,…,0, 1, ,0,…,0, 1,0,…,0, 1,0,…,0,1, 0,…,0,1,0,…,0,1, 2, 0, ..)$

# Feature Extraction

➢ Word Weighting

▪ *Word frequency weighting* and *TF\*IDF weighting*: the number of time that a word appears in a document

▪ Three values are used to calculate the weighting:

- **Term frequency**: the number of time a word appears in a document

- **Collection frequency**: the number of time a word appears in document collection (whole dataset)

- **Document frequency**: the number of document contains a word

=> Features: words have highest Word weighting

# Classifier Evaluation

➤ What is the objective?

- Evaluate quality of the model and its accuracy to know if we can use this model or not

➤ How to do?

- **Accuracy**: the proportion of correctly classified objects
- **Error**: the proportion of incorrectly classified objects
- **Precision:** the proportion of selected items that the system got right
- **Recall:** the proportion of the target items that the system selected
- **Fallout**: the proportion of no targeted items that were mistakenly selected
- **F-measure**: Precision and Recall are combined

# Decision Tree



➤ The first node is root node

➤ Internal nodes are attribute tests

➤ Leaf nodes are class label

➤ Many algorithms ID3, C4.5, CART, CHAID, MARS in decision tree

➤ ID3 uses Entropy and Information Gain

➤ Pruning
  ▪ The pruning step is to avoid **over-fitting**

➤ Cross-validation
  ▪ To maximize the accurate classification of classifier tree model

# Decision Tree (ID3)

➢ Entropy

- Entropy is the indicator of how much information inside a data set

$$Entropy(S) = \sum_{i=1}^{C} -p_i \log_2 p_i$$

Where:

- S: the set of training data
- C: the number of class labels
- $p_i$: the rate of elements belong to class $C_i$

# Decision Tree (ID3)

➢ Information Gain

- Information gain is the measures of reducing entropy in S by an attribute in S

$$Gain(S, A) = Entropy(S) - \sum_{v \in Value(A)} \frac{|S_v|}{S} Entropy(S_v)$$

Where:
- Value (A): set of A values
- $S_v$: subset of S

# Decision Tree (ID3)

➢ General steps of ID3 Algorithm

1. From the dataset S, calculate the entropy of every attribute (feature)
2. Split the set S into subsets using the attribute for which entropy is minimum (or information gain is maximum)
3. Expanding decision tree by adding a node containing that attribute
4. Recursive on subsets using remaining attributes

# Maximum Entropy

➢ Main idea

  ▪ Satisfy constraints

  ▪ Probability distribution of model which is most uniform

➢ What is the constraint?

  ▪ Constraint : If a document contains the word "professor", it has a 40% chance of probability distribution in faculty class

$$f_i(\vec{x}_j, c) = \begin{cases} 1, & if \ w_{ij} > 0 \ and \ c = 1 \\ 0, & otherwise \end{cases}$$

  Wij is the word weighting of word i in document j

# Maximum Entropy

➢ Log-linear Model

  ▪ Use to classify document in Maximum Entropy

$$p(\vec{x}, c) = \frac{1}{Z} \prod_{i=1}^{K} \alpha_i^{f_i(\vec{x}, c)}$$

  Where :

  • K: the number of constraints

  • Z: a constant

  • $\alpha_i$ : the weight of $f_i$

  ▪ Compute $p(\vec{x}_{new}, 1)$ and $p(\vec{x}_{new}, 0)$.

  New document belong to class which has higher probability

# Maximum Entropy

- ➤ Generalized iterative scaling (GIS)
  - ▪ Use to find $\alpha_i$ in the Log-linear Model
  - ▪ GIS find probability distribution which has maximum entropy of Log-linear Model

# Experiment

➢ Dataset : 1000 negative movie reviews and 1000 positive movie reviews

➢ Text Preprocessing

# Experiment

➢ Text representing



Feature and word weighting

# Experiment (result)

```
                worst = 1
                  bring = 0
                    tom = 0
                      details = 0: neg
                      details = 1
                      |    -- = 0: pos
                      |    -- = 1: neg
                    tom = 1
                      come = 0: neg
                      come = 1: pos
                  bring = 1
                    see = 0: neg
                    see = 1
                      usually = 0
                        america = 0: pos
                        america = 1: neg
                      usually = 1: neg
            wonderfully = 1
              red = 0: pos
              red = 1: neg
        stupid = 1
          bob = 0
            into = 0
              perfect = 0
                certainly = 0: neg
                certainly = 1
                  - = 0: pos
                  - = 1: neg
              perfect = 1: pos
            into = 1: neg
          bob = 1
            10 = 0: pos
            10 = 1: neg
```

Classifier Model (ID3)

```
Correctly Classified Instances        247           61.75  %
Incorrectly Classified Instances      153           38.25  %
Kappa statistic                         0.235
Mean absolute error                     0.3825
Root mean squared error                 0.6185
Relative absolute error                76.5    %
Root relative squared error           123.6932 %
Total Number of Instances             400

=== Detailed Accuracy By Class ===

              TP Rate   FP Rate   Precision   Recall   F-Measure   ROC Area   Class
               0.61     0.375      0.619      0.61      0.615       0.618      neg
               0.625    0.39       0.616      0.625     0.62        0.618      pos
Weighted Avg.  0.618    0.383      0.618      0.618     0.617       0.618
```

Result of classifying phase
(training data 66%, test 34%)

# Conclusion

➤ Achievements

- Understand the stages of text classification

- Gather two methods of text classification:

  - Decision Tree method

  - Maximum Entropy method

➤ Future works

- Continue researching methods of text classification
- Program decision tree method to classify document

# Reference

- Christopher D.Manning, Hinrich Schutze, **"Foundations of Statistical Natural Language Processing"**
- Kamal Nigam, John Lafferty, Andrew McCallum, **"Using Maximum Entropy for Text Classification**", In IJCAI-99 Workshop on Machine Learning for Information Filtering
- Kostas Fragos, Yannis Maistros, Christos Skourlas, "**A Weighted Maximum Entropy Language Model for Text Classification**", NLUCS 2005: p.55-67
- Tom M. Mitchell, "**Machine learning**", Published by McGraw-Hill, Maidenhead, U.K., International Student Edition, 1997. ISBN: 0-07-115467-1

# Thank You!