

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Text Categorization

- Problem and Motivation

Hui Fang
Department of Electrical and Computer Engineering
University of Delaware

1

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Text Classification

- A standard supervised learning problem
 - Goal: classify a new document into one of the pre-given categories
 - Training data: labeled document examples

2

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Problem Definition

- Input:
 - a document d
 - a fixed set of categories $C = \{c_1, c_2, \dots, c_J\}$
- Output: a predicted category $c \in C$ for the document d

4

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Email Spam Detection

From: "Laraine" <Claribel@garancs.net>
Subject: Original pilules at factory prices
Date: March 26, 2013 12:05:00 PM EDT
To: <ud-aec@UDel.Edu>

USPS - Fast Delivery Shipping 1-4 day
Best quality drugs
Professional packaging
100% guarantee on delivery
Best prices in the market
Discounts for returning customers
FDA approved products
35000+ satisfied -customers

=====

If you can't click on link, please click "no spam" or copy and paste it to address bar

=====

<http://wawuotu.inhealth-pro.ru>

4

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Examples of Text Classification Problem

- Authorship Identification
 - Classify a document as the one written by different authors
- Sentiment analysis
 - Classify a review to positive or negative
- Automatic essay grading
 - Classify a student essay to A,B,C or D.
- Help desk management system
 - Classify an email according to technical staff's specialty such as Mac, Window, etc.
- ...

5

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Classification methods (1): manually created rules

- Rules based on combinations of words or other features
 - In email spam detection,
 - If an email's sender is in a black-list, the email will be classified as spam.
 - If an email's message contains "dollars" and "have been selected", this email will be classified as spam.
 - If rules are carefully identified and refined by experts, the accuracy can be high.
 - But it requires lots of human efforts to build and maintain the rules.

1

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Classification Methods (2): Supervised Machine Learning

- Input:
 - a document d
 - a fixed set of categories $C = \{c_1, c_2, \dots, c_j\}$

A training set of m hand-labeled documents
 $(d_1, c_1), \dots, (d_m, c_m)$

a learned classifier: $f(d) \rightarrow c$

- Output: a predicted category $c \in C$ for the document d

8

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Automatic Classification Methods

- Generative classifiers
 - Naïve Bayes classifier
- Discriminative classifiers
 - Rocchio classification
 - KNN (K nearest neighbors)
 - Support Vector Machines (SVM)

9

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Feature Selection

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Feature selection → Document Representation

$f(\text{I love this movie. It is sweet and full of humors. The dialogue is great. I love the main actress. She is so pretty. ☺ I would recommend it to my friends. I have watched it multiple times, but I would not mind watching it again.}) = c$

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Bag of Words Representation (Unigram)

$f(\begin{array}{|c|c|} \hline \text{I} & 8 \\ \hline \text{movie} & 3 \\ \hline \text{love} & 3 \\ \hline \text{great} & 4 \\ \hline \dots & \dots \\ \hline \end{array}) = c$

UNIVERSITY of DELAWARE

UNIVERSITY of DE

Use only selected words as features

$f(\text{I love this movie. It is sweet and full of humors. The dialogue is great. I love the main actress. She is so pretty. ☺ I would recommend it to my friends. I have watched it multiple times, but I would not mind watching it again.}) = c$

Use only selected words as features

f(

Love	8
Great	2
Recommend	1
Pretty	1
...	...

)=c