# Text Categorization: Evaluation

Part 1

ChengXiang "Cheng" Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign

# Overview

- What is text categorization?

- Why text categorization?

- How to do text categorization?
  - Generative probabilistic models
  - Discriminative approaches

- **How to evaluate categorization results?**

# General Evaluation Methodology

- Have humans to create a test collection where every document is tagged with the desired categories ("ground truth")
- Generate categorization results using a system on the test collection
- Compare the system categorization decisions with the human-made categorization decisions and quantify their similarity (or equivalently difference)
  - The higher the similarity is, the better the results are
  - Similarity can be measured from different perspectives to understand the quality of results in detail (e.g., which category performs better?)
  - In general, different categorization mistakes may have a different cost that inevitably depends on specific applications, but it is okay not to consider such a cost variation for **relative comparison of methods**

# Classification Accuracy (Percentage of Correct Decisions)

|       | $c_1$  | $c_2$  | $c_3$  | ...  | $c_k$  |
|-------|--------|--------|--------|------|--------|
| $d_1$ | y(+)   | y(-)   | n(+)   |      | n(+)   |
| $d_2$ | y(-)   | n(+)   | y(+)   |      | n(+)   |
| $d_3$ | n(+)   | n(+)   | y(+)   |      | n(+)   |
| ...   |        |        |        |      |        |
| $d_N$ | ...    | ...    |        |      |        |

**+/-  human answer**
(+= correct; - =incorrect)
**y/n    system result**
(y=yes; n=no)

$$\text{Classification Accuracy} = \frac{\text{Total number of correct decisions}}{\text{Total number of decisions made}}$$
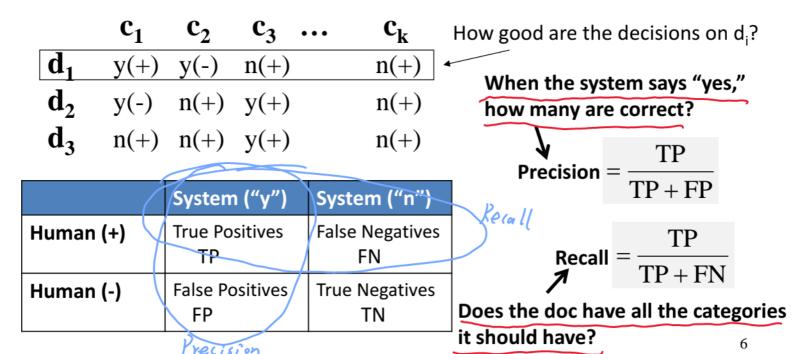
$$= \frac{\text{count}(y(+)) + \text{count}(n(-))}{kN}$$

4

# Problems with Classification Accuracy

- Some decision errors are more serious than others
  - It may be more important to get the decisions right on some documents than others
  - It may be more important to get the decisions right on some categories than others
  - E.g., spam filtering: missing a legitimate email costs more than letting a spam go
- Problem with imbalanced test set
  - Skewed test set: 98% in category 1; 2% in category 2
  - Strong baseline: put all instances in category 1 ➔ 98% accuracy!

# Per-Document Evaluation

|       | $c_1$ | $c_2$ | $c_3$ | $\cdots$ | $c_k$ |
|-------|-------|-------|-------|----------|-------|
| $d_1$ | y(+)  | y(-)  | n(+)  |          | n(+)  |
| $d_2$ | y(-)  | n(+)  | y(+)  |          | n(+)  |
| $d_3$ | n(+)  | n(+)  | y(+)  |          | n(+)  |

How good are the decisions on $d_i$?

**When the system says "yes," how many are correct?**

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

**Does the doc have all the categories it should have?**

|              | System ("y")          | System ("n")           |
|--------------|-----------------------|------------------------|
| **Human (+)**| True Positives TP     | False Negatives FN     |
| **Human (-)**| False Positives FP    | True Negatives TN      |

*Recall*

*Precision*

6

# Per-Category Evaluation

$$c_1 \quad c_2 \quad c_3 \quad \ldots \qquad c_k$$

$d_1$   y(+)   y(-)   n(+)      n(+)

$d_2$   y(-)   n(+)   y(+)      n(+)

$d_3$   n(+)   n(+)   y(+)      n(- )

How good are the decisions on $c_i$?

**When the system says "yes," how many are correct?**

**Precision** $= \dfrac{\text{TP}}{\text{TP} + \text{FP}}$

**Recall** $= \dfrac{\text{TP}}{\text{TP} + \text{FN}}$

|  | System ("y") | System ("n") |
|---|---|---|
| **Human (+)** | True Positives TP | False Negatives FN |
| **Human (-)** | False Positives FP | True Negatives TN |

**Has the category been assigned to all the docs of this category?**

7

# Combine Precision and Recall: F-Measure

$$F_\beta = \cfrac{1}{\cfrac{\beta^2}{\beta^2+1}\cfrac{1}{R}+\cfrac{1}{\beta^2+1}\cfrac{1}{P}} = \cfrac{(\beta^2+1)P*R}{\beta^2 P+R}$$

$$F_1 = \cfrac{2PR}{P+R}$$

**P**: precision
**R**: recall
β: parameter (often set to 1)

Why not 0.5*P+0.5*R?

**What is R if the system says "y" for all category-doc pairs?**

precision 维似低, recall 维很高