**Faculty of Computer Science, Dalhousie University**　　　　*16-Oct-2025*

**CSCI 4152/6509 — Natural Language Processing**

**Lecture 7: Text Classification**

Location: Studley LSC-Psychology P5260　　　Instructor: Vlado Keselj
Time:　　　14:35 – 15:55

**Previous Lecture**

- Collecting n-grams (continued)
- Elements of Information Retrieval
- Vector space model
    - Term weighting schemes:
        - Boolean,
        - *tf* (term frequency, "Bag of Words"),
        - *tf-idf* (term frequency — inverse document frequency)
- Cosine distance measure

## 8.4　Term-by-Document Matrix and Latent Semantic Analysis

Note: This subsection will not be covered in the class.

**Term-by-Document Matrix:**　　The vector space model provides a way to represent each document as a vector. If we have $m$ selected terms for a document collection of $n$ documents, then using for example *tfidf* weights we can represent each of the $n$ documents as an $m$-dimensional vector. If we order these vectors as columns, we get an $m \times n$ dimensional matrix called *term-by-document matrix*. if $w_{ij}$ is the weight of term $t_i$ in the document $d_j$, then the term-by-document matrix is $[w_{ij}]_{m \times n}$, or:

$$
\begin{array}{c|ccccc}
 & d_1 & d_2 & \ldots & d_n \\
\hline
t_1 & w_{11} & w_{12} & \ldots & w_{1n} \\
t_2 & w_{21} & w_{22} & \ldots & w_{2n} \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
t_m & w_{m1} & w_{m2} & \ldots & w_{mn}
\end{array}
\quad \text{or, as a proper matrix:} \quad
\begin{bmatrix}
w_{11} & w_{12} & \ldots & w_{1n} \\
w_{21} & w_{22} & \ldots & w_{2n} \\
\vdots & \vdots & \vdots & \vdots \\
w_{m1} & w_{m2} & \ldots & w_{mn}
\end{bmatrix}
$$

**Dimensionality reduction:**　　The number of different terms generally corresponds to the number of different words in a document collection, and this number is generally large, in the range of 100,000s. It is useful for various application to reduce this dimensionality. Some ways of reducing dimensionality are: removing stop-words and very rare words, and selecting only the most distinctive terms, which is a process known as feature selection.

**Latent Semantic Analysis:** An interesting mathematical way of representing documents in a vector space with much lower dimensionality is known as Latent Semantic Analysis.

**Latent Semantic Analysis**

- A method for term-by-document dimensionality reduction
- Also known as Latent Semantic Indexing (LSI) in IR

  – Example with four terms and two documents
  – Main idea: use singular value decomposition on term-by-document matrix $M$
  – Singular value decomposition: $M_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}^T$
  – Closest by Frobenius norm matrix of rank $\leq k$ is
    $M_{m \times n}^{(k)} = U_{m \times m} \Sigma_{m \times n}^{(k)} V_{n \times n}^T$
  – Concept and document representations

## 8.5   IR Evaluation Measures: Precision, Recall, and F-measure

Note: This section is normally covered in an earlier Machine Learning course, so it is covered in the class as a review material.

We will now define some main evaluation measures used in IR, which are also important in general text mining tasks, such as text classification. The main three measures are: precision, recall, and F-measure:

  – **Precision** is the percentage of true positives out of all returned documents; i.e.,

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

  – **Recall** is the percentage of true positives out of all relevant documents in the collection; i.e.,

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

  – **F-measure** is a weighted harmonic mean between Precision and Recall:

$$F = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

  – We usually set $\beta = 1$, in which case we have:

$$F = \frac{2PR}{P + R}$$

Precision and Recall can be explained in the following way: Given a query, the search engine identifies a set of relevant documents, and returns this set. Some of the returned documents are truly relevant and we call them *true positives (TP);* some returned documents are not relevant and we call them *false positives (FP);* some documents are relevant but were not returned by the engine and we call them *false negatives (FN);* and the remaining documents are not returned by the engine and they are not relevant, and we call them *true negatives (TN).*

The typical value in the F-measure is $\beta = 1$, for equal emphasis on precision and recall. However, if we want to put more emphasis on precision we choose $\beta$ close to 0, and for more emphasis on recall we choose $\beta$ close to 1. $\beta$ must always be from the interval $[0, 1]$, i.e., $0 \leq \beta \leq 1$.
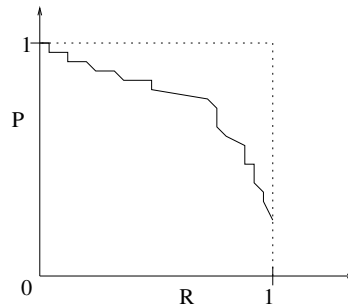
## 8.6   Recall-Precision Curve

Note: Covered in class as review material.

Precision, recall and F-measure may be seen as too simplistic views of evaluating a search engine since a search engine does not return just a set of relevant documents, but a ranked list of relevant documents.

A more appropriate way to evaluate this ranked list is using the *recall-precision curve.* The basic idea of the recall-precision curve is to draw a point in a two-dimensional plane corresponding to precision ($y$ axis) and recall ($x$ axis) of the following document sets: first ranked document, the first two ranked documents, the first three ranked documents, and so on. Although such curve is a roughly smooth curve going from the point (0,1) to (1,0), it does usually have some noisy changes of direction, so it is "smoothed" by actually using the interpolated precision curve.

A typical Recall-Precision curve looks as follows:



To avoid "noisy" changes of curve direction, *interpolated precision* (*IntPrec*) is often used. Interpolated precision is the maximal precision that is obtained with certain recall level; namely, if $P(k)$ and $R(k)$ are precision and recall of the set of first $k$ ranked documents, than for any recall value $r \in [0, 1]$:

$$IntPrec(r) = \max_{k, R(k) \geq r} P(k)$$

**Recall-Precision Curve Example.** Suppose that a search engine returned 12 ranked results to our query, and when we checked them, the following are our judgments on their relevance:

1. relevant
2. relevant
3. relevant
4. not relevant
5. relevant
6. not relevant
7. relevant
8. not relevant
9. not relevant
10. relevant
11. not relevant
12. not relevant

**Task 1: Precision, Recall and F-measure**

 – Assuming that the total number of relevant documents in the collection is 8, calculate precision, recall, and F-measure ($\beta = 1$) for the returned 12 results.

Since there is a total of 6 relevant documents among the set of 12, we can calculate precision to be $P = \frac{6}{12} = 0.5$.

It is assumed that there is a total number of 8 relevant documents in the collection, so the recall is $R = \frac{6}{8} = 0.75$.

Finally, we calculate the F-measure: $F = \frac{2PR}{P+R} = \frac{2 \cdot 0.5 \cdot 0.75}{0.5 + 0.75} = \frac{0.75}{1.25} = \frac{3}{5} = 0.6$.

**Task 2: Recall-Precision Curve**

 – Task: Draw the recall-precision curve for these results
 – First step: Form sets of $n$ initial documents, and look at their relevance:

Set 1: $\{R\}$:                                    $R_1 = \frac{1}{8} = 0.125$    $P_1 = \frac{1}{1} = 1$
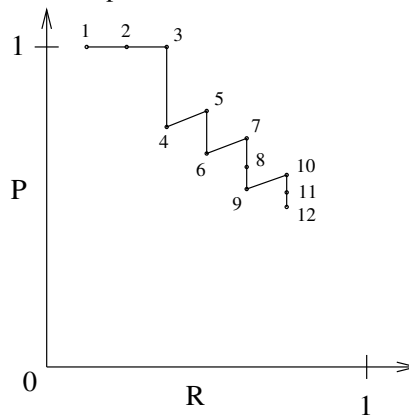
Set 2: $\{R, R\}$:                                 $R_2 = \frac{2}{8} = 0.25$    $P_2 = \frac{2}{2} = 1$

Set 3: $\{R, R, R\}$:                              $R_3 = \frac{3}{8} = 0.375$    $P_3 = \frac{3}{3} = 1$

Set 4: $\{R, R, R, NR\}$:                                      $R_4 = \frac{3}{8} = 0.375$   $P_4 = \frac{3}{4} = 0.75$

Set 5: $\{R, R, R, NR, R\}$:                                   $R_5 = \frac{4}{8} = 0.5$   $P_5 = \frac{4}{5} = 0.8$

Set 6: $\{R, R, R, NR, R, NR\}$:                               $R_6 = \frac{4}{8} = 0.5$   $P_6 = \frac{4}{6} \approx 0.666666666666667$

Set 7: $\{R, R, R, NR, R, NR, R\}$:                            $R_7 = \frac{5}{8} = 0.625$   $P_7 = \frac{5}{7} \approx 0.714285714285714$

Set 8: $\{R, R, R, NR, R, NR, R, NR\}$:                        $R_8 = \frac{5}{8} = 0.625$   $P_8 = \frac{5}{8} = 0.625$

Set 9: $\{R, R, R, NR, R, NR, R, NR, NR\}$:                    $R_9 = \frac{5}{8} = 0.625$   $P_9 = \frac{5}{9} \approx 0.555555555555556$

Set 10: $\{R, R, R, NR, R, NR, R, NR, NR, R\}$:                $R_{10} = \frac{6}{8} = 0.75$   $P_{10} = \frac{6}{10} = 0.6$

Set 11: $\{R, R, R, NR, R, NR, R, NR, NR, R, NR\}$:            $R_{11} = \frac{6}{8} = 0.75$   $P_{11} = \frac{6}{11} \approx 0.545454545454545$

Set 12: $\{R, R, R, NR, R, NR, R, NR, NR, R, NR, NR\}$:        $R_{12} = \frac{6}{8} = 0.75$   $P_{12} = \frac{6}{12} = 0.5$

Using these ten points, we can draw the recall-precision curve:



The recall-precision curve that we just saw is not exactly monotonically non-increasing although the general trend of these curves is to grom from near the point $(0, 1)$ down towards point $(1, 0)$ in the coordinate system. The reason is that after the initial 100% precision, each time we see a new relevant document in the ranked list both precision and recall will increase a bit, and when we see a non-relevant document the precision will drop with the same recall, which creates a bit of a zig-zagged line. To make the non-increasing, and thus a bit smoother, we use the *interpolated precision-recall curve,* which is obtained by keeping precision at the level that is maximal from the associated recall point and forward.

*Slide notes:*

**Task 3: Interpolated Recall-Precision Curve**
  – Task: Draw interpolated Recall-Precision curve
  – Formula:
$$IntPrec(r) = \max_{k, R(k) \geq r} P(k)$$

  – Based on the previous Task:
    $0 \leq r \leq R_4 = \frac{3}{8} = 0.375 \Rightarrow IntPrec(r) = 1$
    $R_4 < r \leq R_6 = \frac{4}{8} = 0.5 \Rightarrow IntPrec(r) = 0.8$
    $R_6 < r \leq R_9 = \frac{5}{8} = 0.625 \Rightarrow IntPrec(r) = 5/7 \approx$
    $0.714285714$
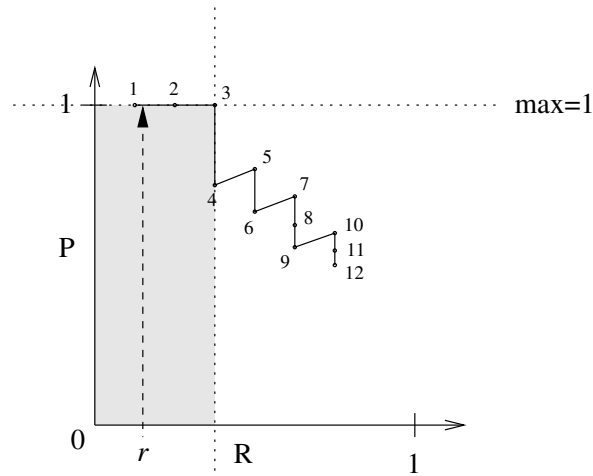    $R_9 < r \leq R_{12} = \frac{6}{8} = 0.75 \Rightarrow IntPrec(r) = 0.6$

To calculate the interpolated recall-precision curve, we use the formula:

$$IntPrec(r) = \max_{x, R(k) \geq r} P(k)$$

To use this formula, we can start first with $r = 0$, which gives:

$$IntPrec(0) = \max_{k, R(k) \geq 0} P(k) = \max_k P(k) = 1$$

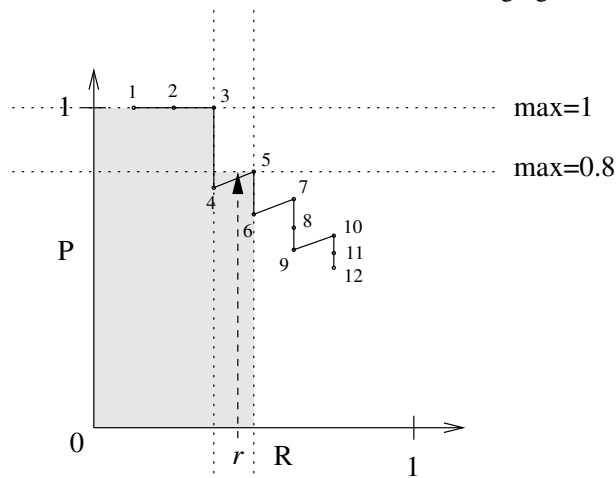because $R(k) \geq 0$ for all $k$, and maximum $P(k)$ is 1. If we increase $r$ starting from 0, we see that the maximal precision will remain 1 for all points $R_1$, $R_2$, $R_3$, and $R_4$, as shown in the following figure:



This is how we get the following values for the Interpolated Precision:
$0 \leq r \leq R_4 = \frac{3}{8} = 0.375 \Rightarrow IntPrec(r) = 1$

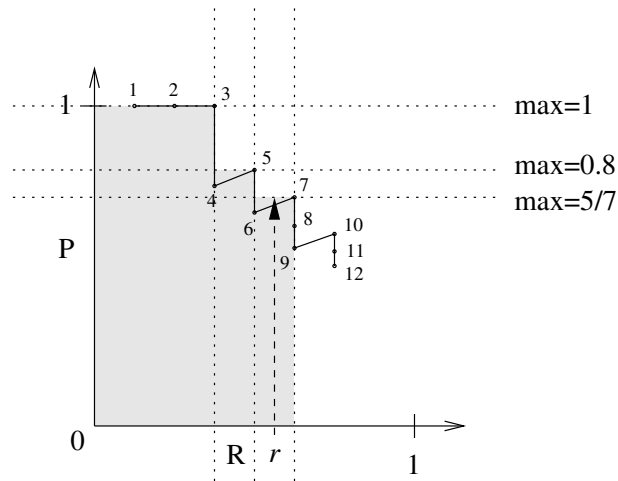For the values $r > R_4$, the next maximum is 0.8, as shown in the following figure:



and that is how we obtain the next interval:

$R_4 < r \leq R_6 = \frac{4}{8} = 0.5 \Rightarrow IntPrec(r) = 0.8$

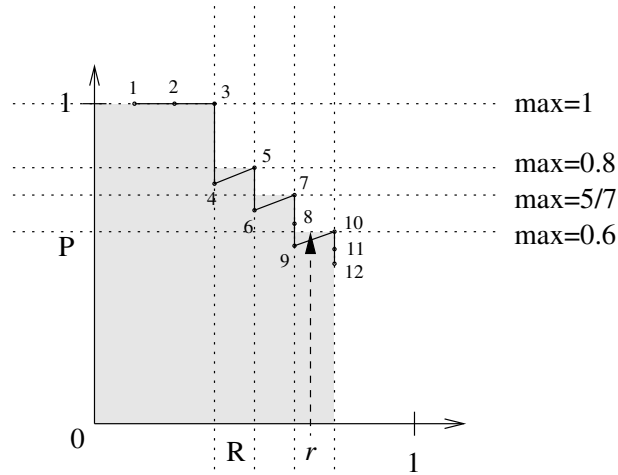Similarly, for the values $r > R_6$, the next maximum is $5/7$, as shown in the following figure:

and that is how we obtain the next interval:

$R_6 < r \le R_9 = \frac{5}{8} = 0.625 \Rightarrow IntPrec(r) = 5/7 \approx 0.714285714$

For the values $r > R_9$, the next maximum is 0.6, as shown in the following figure:



and that is how we obtain the final interval:

$R_9 < r \le R_{12} = \frac{6}{8} = 0.75 \Rightarrow IntPrec(r) = 0.6$

There are no further points, so we can finish the curve at this point. We can summarize the values for the interpolated recall-precision curve as:
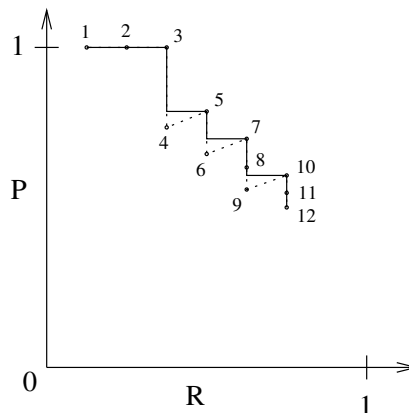
$0 \le r \le R_4 = \frac{3}{8} = 0.375 \Rightarrow IntPrec(r) = 1$
$R_4 < r \le R_6 = \frac{4}{8} = 0.5 \Rightarrow IntPrec(r) = 0.8$
$R_6 < r \le R_9 = \frac{5}{8} = 0.625 \Rightarrow IntPrec(r) = 5/7 \approx 0.714285714$
$R_9 < r \le R_{12} = \frac{6}{8} = 0.75 \Rightarrow IntPrec(r) = 0.6$

Using these points, we can construct the interpolated precision curve, up to $R = 0.75$:

A way to look at the interpolated recall-precision curve is that it is obtained from the direct recall-precision curve by "filling up" all parts where the curve increases, in a way similar as water would fill up the "depressions" if we imagine the curve to be a waterfall cascade.

**Interpolated R-P Curve at 11 Standard Levels**

**Some Other Similar Measures**

    – Fallout

$$Fallout = \frac{FP}{FP + TN}$$

    – Specificity

$$Specificity = \frac{TN}{TN + FP}$$

    – Sensitivity

$$Sensitivity = \frac{TP}{TP + FN} \quad (= R)$$

    – Sensitivity and Specificity: useful in classification and contexts such as medical tests

Sensitivity and Specificity are typically used in the classification task, that we will described later. They are also frequently useful to evaluate medical tests. For example, if we consider a context of a medical test for disease diagnostics: a sensitive test, same as recall, is good at not missing any true cases of a disease (true positives); while a specific test is good in eliminating a possible suspected disease.

# 9 Text Classification as General NLP Task

Note: Covered in class as review material.

## 9.1 Text Classification as a Text Mining Task

Text Classification is one of the tasks in a more general area called *Text Mining*. The area of Text Mining generally deals with processing of large quantities of text and deriving some useful information, knowledge, or insight from it. The name Text Mining is derived from a similar area of Data Mining, which deals with large quantities of data in general for the purpose of knowledge discovery. Some tasks in Text Mining are related to similar tasks in a wider Data Mining area. The following are some typical text mining tasks:

- Text Classification
- Text Clustering
- Information Extraction
- And some new and less prominent tasks:
    - Text Visualization
    - Filtering tasks, Event Detection
    - Terminology Extraction

*Text Classification* is the task of classifying documents into classes of documents; i.e., sets of documents, of certain properties. For example, classifying email into spam email and non-spam email is an example of text classification. *Text Clustering* is the task of grouping documents in a collection in a groups of similar documents, or clusters. *Information Extraction* is the task of extracting table-like data from text documents. For example, processing news and filling out information about companies, their names, addresses, and names of CEOs would be an information extraction task.

The area of text visualization addresses the problem of different visual representations of text and textual documents. Filtering tasks deal with selecting relevant documents or information from a stream of usually short textual documents. Event detection is the task of detecting events from a stream of documents, such as news-wire. A CEO change in a company could be one kind of interesting events to be detected. Terminology Extraction is the task of extracting terminology, i.e., domain terms, from a document collection in a domain. For example, analyzing bio-medical scientific papers and detecting and extracting terminology such as protein names would be an example of terminology extraction.

## 9.2   Types of Text Classification

**Text Classification** is also known as Text Categorization. It is the problem of automatically classifying a document into one of predetermined classes or categories of documents. In a more usual form of classification, we always assign a document to exactly one class. In a more flexible form, known as *multi-label* classification, we assign document to zero or more classes; or we can view the task as assigning a set of labels to the document, where each label is a designation of a class.

Beside some reading on document classification in the textbook [JM], there is a more elaborate description in the Manning and Schütze book ([MS]), in Chapter 16: Text Categorization.

**Types of Text Classification**

- topic categorization
- sentiment classification
- authorship attribution and plagiarism detection
- authorship profiling (e.g., age and gender detection)
- spam detection and e-mail classification
- encoding and language identification
- automatic essay grading

More specialized example: dementia detection using spontaneous speech

**Creating Text Classifiers**

- Can be created manually
    - typically rule-based classifier
    - example: detect or count occurrences of some words, phrases, or strings
- Another approach: make programs that *learn* to classify
    - In other words, classifiers are generated based on labeled data

– supervised learning

While we can create a classifier from the scratch, a more usual approach is to classify, i.e., label, a set of document manually, and then devise a method generate a classifier based on this set of labeled documents. This process of generating a classifier is known as *training*, or *machine learning*. The classification problem is known as an example of *supervised learning* in the machine learning area, since we need to provided labeled examples; in other words, we "supervise" the learning algorithm.

## 9.3  Evaluation Measures for Text Classification

When we build a text classifier, an important question is how to evaluate it so we can measure how good it is and how different classifiers compare. In order to do this, we first need to prepare a dataset consisting of documents, where each document is labeled with the class that it belongs too. This is called a *labeled dataset* or a *labeled document collection.* The labels are usually assigned and checked manually, and we take them as the *ground truth,* also called the *gold standard,* against which we evaluate the classifier.

We run the classifier against the documents, with labels being hidden, and once the classifier assigns labels to the documents, we compare them against the gold standard labels. The first evaluation measure that we can calculate is the *accuracy,* which is the percentage or fraction of the documents being correctly classified by the classifier. If we want to examine in more details how classifier performed on documents from different classes, we present the results in the *contingency table* also known as the *confusion matrix.*

*Slide notes:*

**Evaluation Measures for Text Classification**

– Contingency table (confusion matrix) and Accuracy
– Example (classes $A$, $B$, and $C$):

|  |  | Gold standard | | | |
|---|---|---|---|---|---|
|  |  | $A$ | $B$ | $C$ | |
| Model | $A$ | 5 | 1 | 1 | 7 |
| classification | $B$ | 3 | 10 | 2 | 15 |
|  | $C$ | 0 | 2 | 10 | 12 |
|  |  | 8 | 13 | 13 | 34 |

– Accuracy: percentage of correct classifications; in the example, $= 25/34 \approx 0.7353 = 73.53\%$

**Per class: Precision, Recall, and F-measure**

– For each class: Yes = in class, No = not in class

|  | Yes is correct | No is correct |
|---|---|---|
| Yes assigned | $a$ | $b$ |
| No assigned | $c$ | $d$ |

– precision ($\frac{a}{a+b}$), recall ($\frac{a}{a+c}$), fallout ($\frac{b}{b+d}$), F-measure:

$$F = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

– If $\beta = 1 \Rightarrow$ Precision and Recall treated equally
– macro-averaging (equal weight to each class) and micro-averaging (equal weight to each object) (2×2 contingency tables vs. one large contingency table)

**Example:** Let us assume that we are evaluating an authorship-attribution classifier. The classifier is presented with 34 documents, written by three authors: A1, A2, and A3, and the classifier produces labels. We know the

true authorship of the documents, so we compare the classifier results to these labels. The obtained results can be presented as a so-called *confusion matrix*, or *contingency table*:

|  |  | Gold standard | | | |
|---|---|---|---|---|---|
|  |  | A1 | A2 | A3 |  |
| System | A1 | 5 | 1 | 1 | 7 |
| response | A2 | 3 | 10 | 2 | 15 |
|  | A3 | 0 | 2 | 10 | 12 |
|  |  | 8 | 13 | 13 | 34 |

Or, we can create contingency tables for each class separately:

|  | Gold standard | | |
|---|---|---|---|
|  | A1 | not A1 |  |
| A1 | 5 | 2 | 7 |
| not A1 | 3 | 24 | 27 |
|  | 8 | 26 | 34 |

|  | Gold standard | | |
|---|---|---|---|
|  | A2 | not A2 |  |
| A2 | 10 | 5 | 15 |
| not A2 | 3 | 16 | 19 |
|  | 13 | 21 | 34 |

|  | Gold standard | | |
|---|---|---|---|
|  | A3 | not A3 |  |
| A3 | 10 | 2 | 12 |
| not A3 | 3 | 19 | 22 |
|  | 13 | 21 | 34 |

The overall accuracy can be calculated using the overall table;

$$Accuracy = \frac{5 + 10 + 10}{34}$$

Per-class precisions are:

$$P_{A1} = \frac{5}{7} \quad P_{A2} = \frac{10}{15} \quad P_{A3} = \frac{10}{12}$$

Per-class recalls are:

$$R_{A1} = \frac{5}{8} \quad R_{A2} = \frac{10}{13} \quad R_{A3} = \frac{10}{13}$$

Macro-averaged precision, recall, and F-measure are:

$$P_{macro} = \frac{5/7 + 10/15 + 10/12}{3} \quad R_{macro} = \frac{5/8 + 10/13 + 10/13}{3} \quad F_{macro} = \frac{2 \cdot P_{macro} \cdot R_{macro}}{P_{macro} + R_{macro}}$$

To calculate micro-averaged precision, recall, and F-measure, we calculate cumulative per-class table:

|  | Gold standard | | |
|---|---|---|---|
|  | A | not A |  |
| A | 25 | 9 | 34 |
| not A | 9 | 59 | 68 |
|  | 34 | 68 | 102 |

and then we calculate the micro-averaged measures:

$$P_{micro} = \frac{25}{34} \quad R_{micro} = \frac{25}{34} \quad F_{micro} = \frac{2 \cdot P_{micro} \cdot R_{micro}}{P_{micro} + R_{micro}} = \frac{25}{34}$$