

# On the Use of Reliable-Negatives Selection Strategies in the PU Learning Approach for Quality Flaws Prediction in Wikipedia

---

Edgardo Ferretti and Marcelo Errecalde

Universidad Nacional de San Luis  
{ferretti,merreca}@unsl.edu.ar

Maik Anderka

University of Paderborn  
maik.anderka@uni-paderborn.de

Benno Stein

Bauhaus-Universität Weimar  
benno.stein@uni-weimar.de

11<sup>th</sup> International Workshop on Text-based Information Retrieval, TIR'14  
Munich, Germany

September 4th, 2014

# Information Quality in Wikipedia

## Situation

- ❑ extremely varying content quality
  - everyone can edit Wikipedia, even anonymously
  - heterogeneous community of Wikipedia authors
  - edits are not reviewed before publication
- ❑ comprehensive manual quality assurance is unfeasible
  - large data volumes, constantly evolving contents



# Information Quality in Wikipedia



## Situation

- ❑ extremely varying content quality
  - everyone can edit Wikipedia, even anonymously
  - heterogeneous community of Wikipedia authors
  - edits are not reviewed before publication
- ❑ comprehensive manual quality assurance is unfeasible
  - large data volumes, constantly evolving contents

## Previous work

- ❑ research question: “Is an article featured or not?”  
[Hu et al., CIKM’07] [Blumenstock, WWW’08] [Dalip et al., JCDL’09] [Lipka and Stein, WWW’10]
- no practical support for Wikipedia’s quality assurance process
- less than 0.1% of the English Wikipedia articles are featured

# Quality Flaw Prediction in Wikipedia

## Question

- ❑ How to improve the 99.9% non-featured Wikipedia articles?

## Central idea

- ❑ automatic exploitation of human-defined cleanup tags [Anderka et al., WWW'11]

The screenshot shows the Wikipedia article for "BASE jumping". A prominent orange-bordered warning box at the top of the article text states: "This article does not cite any references or sources. Please help improve this article by adding citations to reliable sources. Unsourced material may be challenged and removed. (February 2010)". The article text includes several instances of "[citation needed]" tags, which are also circled in orange. These tags appear in the "History" section and in the main text: "judging criteria", "Recent years have seen a formal competition held at the 452 metres (1,483 ft) high Petronas Towers in Kuala Lumpur, Malaysia, judged on landing accuracy", and "In 2010 Northern Norway celebrated with a world record with 53 Base jumpers jumping from a cliff". The article also features a sidebar with navigation links, a search bar, and a small image of a person base jumping from a cliff.

# Quality Flaw Prediction in Wikipedia

## Question

- ❑ How to improve the 99.9% non-featured Wikipedia articles?

## Central idea

- ❑ automatic exploitation of human-defined cleanup tags [Anderka et al., WWW'11]
  - each tag defines a specific quality flaw
  - tagged articles serve as human-labeled examples
  - machine learning is used to predict flaws in untagged articles

## Existing flaw prediction approaches

- ❑ one-class classification [Anderka et al., WWW'11, SIGIR'12]
- ❑ binary classification [Ferschke et al., CLEF'12, ACL'13]
- ❑ **PU learning** [Ferretti et al., CLEF'12]

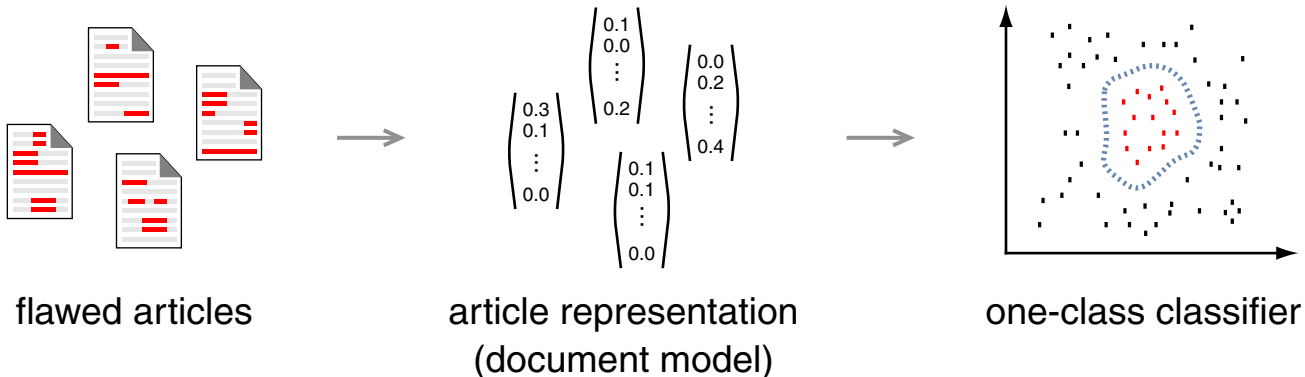
# Outline

- Motivation
- Problem Statement
- Quality Flaw Prediction Using PU Learning
- Analysis and Empirical Evaluation
- Summary

# Problem Statement

## Quality flaw prediction in Wikipedia [Anderka et al., SIGIR'12]

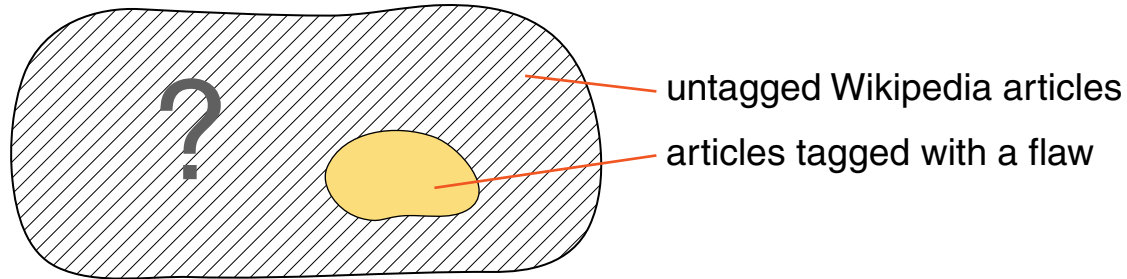
- 3.8 M English Wikipedia articles  $\rightarrow D$
- 445 quality flaws (cleanup tags)  $\rightarrow F$
- Build a classifier  $c : D \rightarrow \{1; 0\}$  for each flaw  $f \in F$ , given a sample of articles containing  $f$ .



# Problem Statement

## Quality flaw prediction using PU learning [Ferretti et al., CLEF'12]

- exploit untagged articles to improve the effectiveness of a classifier  $c$



- in Wikipedia, it is more than likely that many flaws are not yet identified
- PU learning: learning from *Positive* and *Unlabeled* examples [Liu et al., ICML'02]
- *positive* examples = articles tagged with a flaw
  - *unlabeled* examples = untagged articles (either flawed or flawless)



# Problem Statement

Background: PU learning [Liu et al., ICML'02]

- set  $P$  of positive examples
- set  $U$  of unlabeled examples (containing both positive and negative examples)
- Build a classifier using  $P$  and  $U$  that can identify positive examples in  $U$  or in a separate test set.
  
- two-stage approach:
  1. identifying *reliable negatives*
    - train a binary classifier using  $P$  and  $U$
    - apply this classifier to the examples in  $U$
    - consider all examples not classified as “positive” as *reliable negatives*
  
  2. building the final classifier (non-iterative version)
    - train a binary classifier using  $P$  and the set of *reliable negatives*

# Problem Statement

## Crucial aspects in the Wikipedia setting

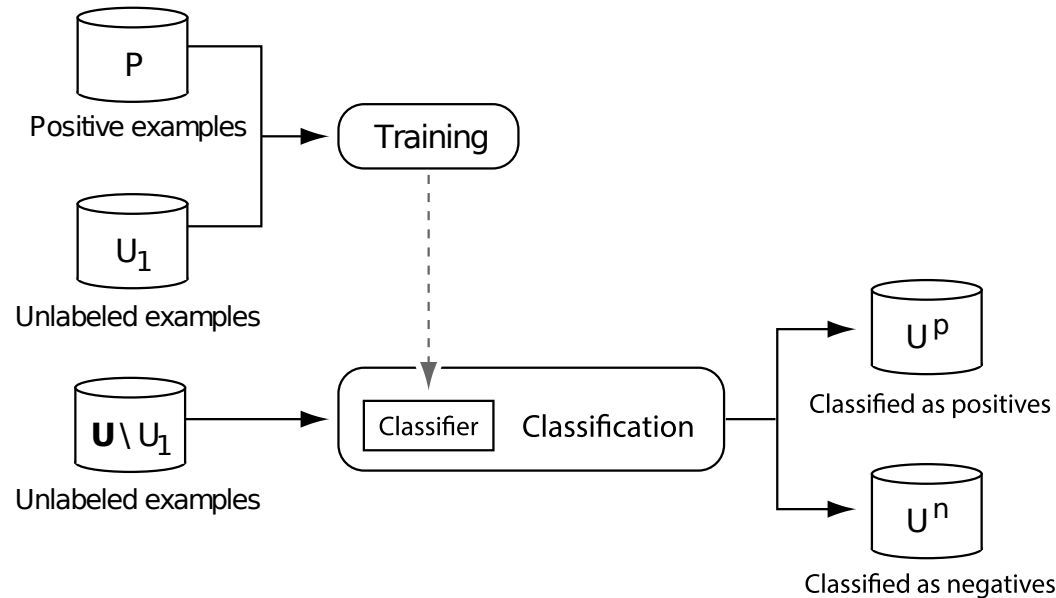
1. unknown (flaw-specific) class imbalances
    - 1<sup>st</sup> stage: ratio between  $P$  and  $U$
    - 2<sup>nd</sup> stage: ratio between  $P$  and the set of *reliable negatives*
  2. effects of sampling (essential in practice due to the large number of existing Wikipedia articles)
    - 1<sup>st</sup> stage:  $U$  is very large for most flaws
    - 2<sup>nd</sup> stage: the set of *reliable negatives* can become considerably large
- have not—or only partially—addressed by Liu et al. and Ferretti et al.
  - we show where in the PU learning procedure sampling is useful
  - we analyze how different sampling strategies affect the flaw prediction effectiveness

# Outline

- Motivation
- Problem Statement
- Quality Flaw Prediction Using PU Learning
- Analysis and Empirical Evaluation
- Summary

# Quality flaw prediction using PU learning

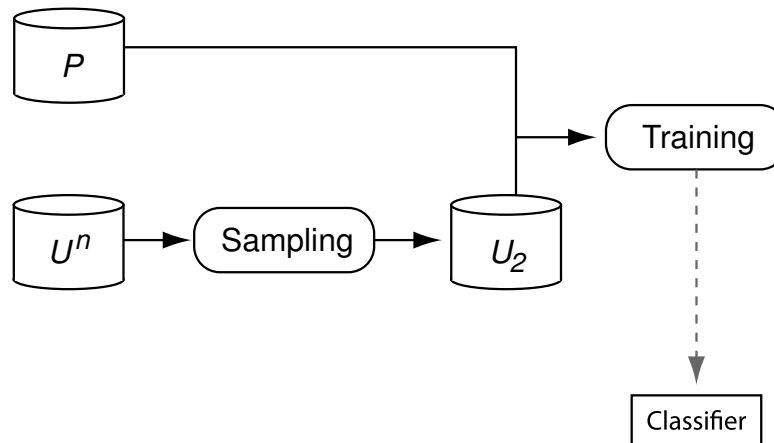
1<sup>st</sup> stage: identifying *reliable negatives*



- $U_1$  is a sample from  $U$
- training set is balanced,  $|P| = |U_1|$
- sampling strategy does not affect the flaw prediction performance
- random sampling

# Quality flaw prediction using PU learning

2<sup>st</sup> stage: building the final classifier



□ using  $U_2 = U^n$  worsened the performance by up to 50% [Ferretti et al., CLEF'12]

□ sampling strategies:

$M_1$  selecting  $|P|$  articles by random from  $U^n$

$M_2$  selecting the  $|P|$  *best* articles from  $U^n$

(those assigned the highest confidence values by the first-stage classifier)

$M_3$  selecting the  $|P|$  *worst* articles from  $U^n$

(those assigned the lowest confidence values by the first-stage classifier)

# Outline

- Motivation
- Problem Statement
- Quality Flaw Prediction Using PU Learning
- Analysis and Empirical Evaluation
- Summary

# Analysis and Empirical Evaluation

## Experimental design

- evaluation corpus of the “1<sup>st</sup> international competition on quality flaw prediction in Wikipedia”
  - 1,592,226 English Wikipedia articles
  - 208,228 tagged to contain one of ten important quality flaws
- 1<sup>st</sup> stage classifier: Naïve Bayes
- 2<sup>nd</sup> stage classifier: Support Vector Machine (SVM)
- balanced training sets:  $|P| = |U_1|$  and  $|P| = |U_2|$
- random sampling in the 1<sup>st</sup> stage
- $M_1$ ,  $M_2$ , and  $M_3$  in the 2<sup>nd</sup> stage

# Analysis and Empirical Evaluation

Selecting *reliable negatives* (2<sup>nd</sup> stage sampling)

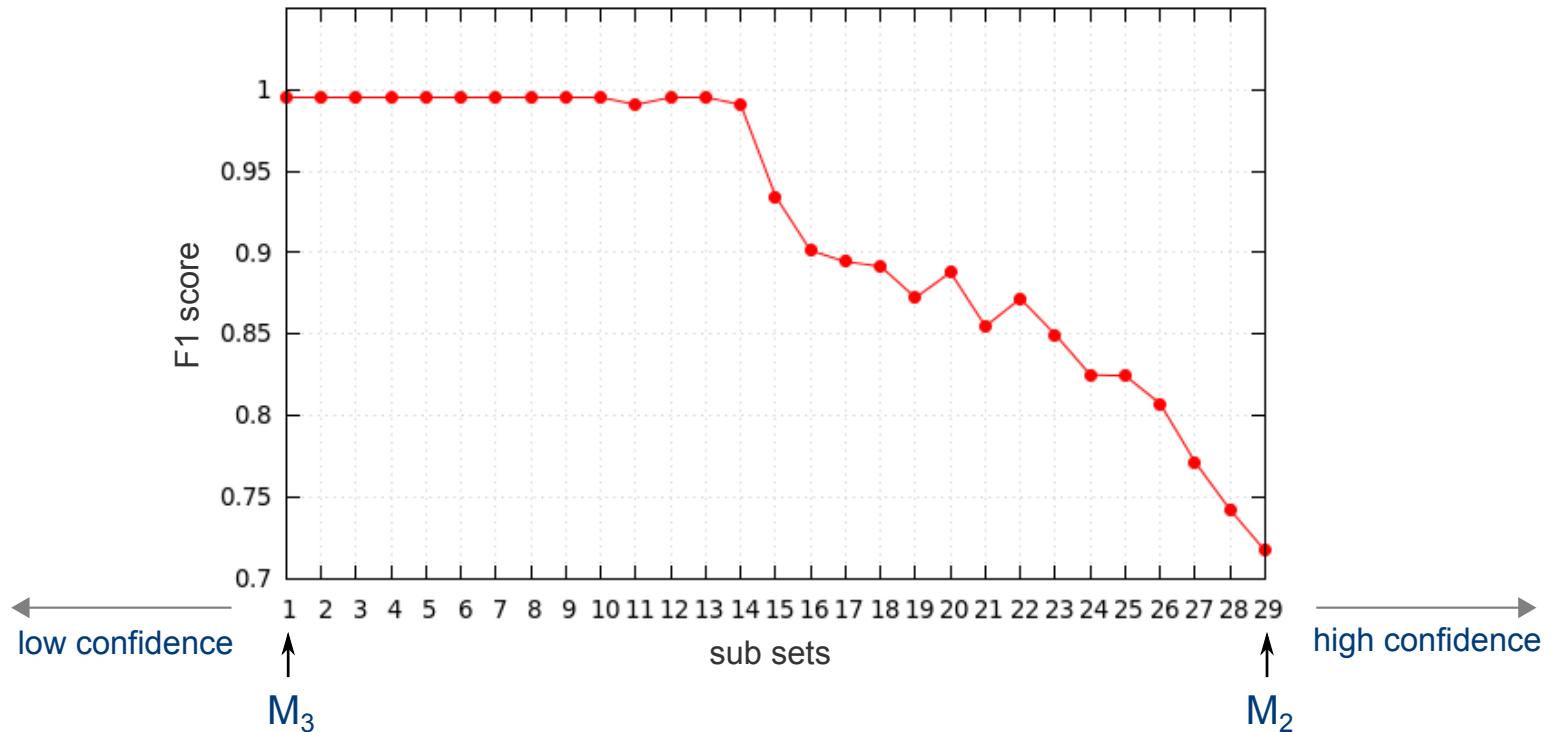
- flaw *Unreferenced*:  $|U^n| = 29,635$ ,  $|P| = |U_2| = 1,000$



# Analysis and Empirical Evaluation

## Selecting *reliable negatives* ( $2^{\text{nd}}$ stage sampling)

- flaw *Unreferenced*:  $|U^n| = 29,635$ ,  $|P| = |U_2| = 1,000$



→ strategy  $M_3$  outperforms  $M_2$

→ differences between  $M_3$  and  $M_1$  (random) are not statistically significant

# Analysis and Empirical Evaluation

## Flaw prediction effectiveness

effectiveness of PU learning in terms of F1 score for the ten quality flaws

<b>flaw name</b>	<b>baseline</b> [Ferretti et al., CLEF'12]	<b>proposed approach</b> using strategy $M_3$
<i>Advert</i>	0.8214	0.9440 (+14.93%)
<i>Empty section</i>	0.8216	0.9394 (+14.34%)
<i>No footnotes</i>	0.8264	0.9826 (+18.90%)
<i>Notability</i>	0.7944	0.9886 (+24.45%)
<i>Orphan</i>	0.8986	0.9960 (+10.84%)
<i>Original research</i>	0.7638	0.9338 (+22.26%)
<i>Primary sources</i>	0.8068	0.9891 (+22.60%)
<i>Refimprove</i>	0.8362	0.9382 (+12.20%)
<i>Unreferenced</i>	0.8365	0.9432 (+12.76%)
<i>Wikify</i>	0.7396	0.9818 (+32.75%)
<b>averaged over all flaws</b>	<b>0.8145</b>	<b>0.9637 (+18.31%)</b>

# Outline

- ❑ Motivation
- ❑ Problem Statement
- ❑ Quality Flaw Prediction Using PU Learning
- ❑ Analysis and Empirical Evaluation
- ❑ **Summary**

# Summary

## What we have done

1. shed light on the effects of sampling in PU learning
  - sampling is necessary (in both stages)
  - in general, sampling strategy  $M_3$  is favorable
2. improved PU learning approach for quality flaw prediction in Wikipedia
  - average improvement of 18.31% compared to the baseline

# Summary

## What we have done

1. shed light on the effects of sampling in PU learning
  - sampling is necessary (in both stages)
  - in general, sampling strategy  $M_3$  is favorable
2. improved PU learning approach for quality flaw prediction in Wikipedia
  - average improvement of 18.31% compared to the baseline

## Current work

- comparative study of the existing flaw prediction approaches

# Thank you!

[maik.anderka@uni-paderborn.de](mailto:maik.anderka@uni-paderborn.de)

# Appendix

# Article representation

- 65 state-of-the-art features, 30 new features

**content** characters, words, syllables, sentences, readability, parts of speech, closed-class word sets, . . .

**structure** sections, tables, images, references, categories, templates, lists, specific sections, . . .

**network** internal-, external-, interwiki-, broken links, PageRank, citation measures, . . .

**edit history** age, currency, connectivity, revisions, reverts, editors, cooperation, . . .