# Measuring Quality of Collaboratively Edited Documents :

## the case of Wikipedia

# The presentation is given by :

**Amin El Iraki**

**Natixis**

———————

**Data Scientist**

———————

@ amine.el-iraki@dauphine.eu

**Olivier Randavel**

**M13H**
**(Start-up) Data Marketing**

———————

**Data Scientist**

———————

@ olivier.randavel@gmail.com

**Hadrien Mariaccia**

**Dassault Systèmes**

———————

**Data Scientist**

———————

@ hadrien.mariaccia@dauphine.eu

**Hippolyte Mayard**

**Société Générale**

———————

**Data Scientist**

———————

@ hippolyte.mayard@dauphine.eu

**Valentin Vu Van**

**Thales**

———————

**Data Scientist**

———————

@ valentin.vu-van@dauphine.eu

# Measuring Quality of Collaboratively Edited Documents :
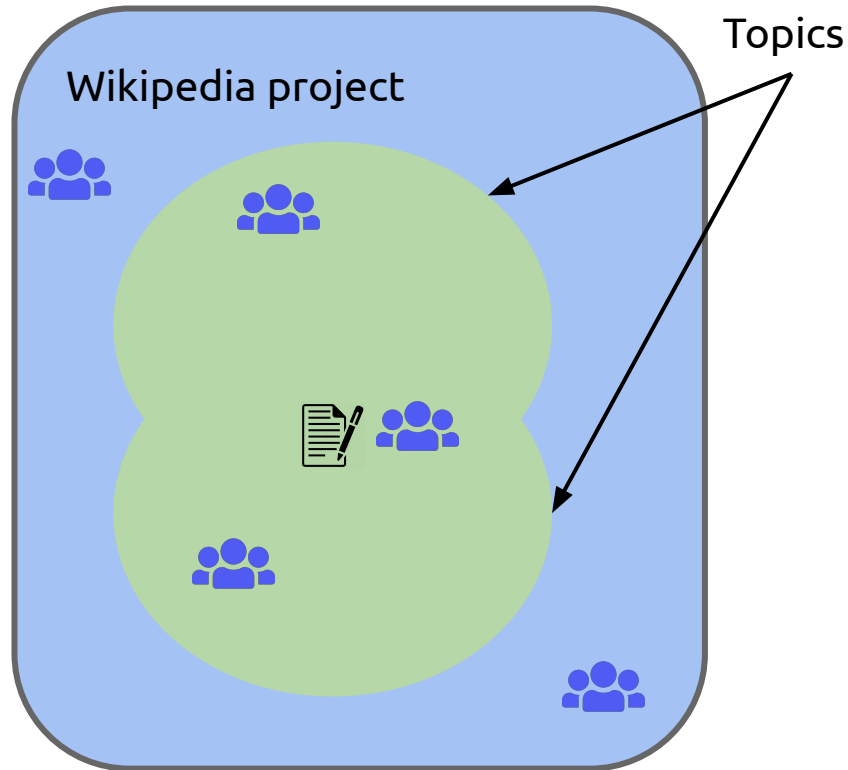
## the case of Wikipedia

# Use case

# Table of contents

1. Wikipedia ranking model

1. Research paper model

1. Pros and cons

# Wikipedia presentation :

Topics

**Wikipedia project**

49 millions articles, 6 millions in English

+300 Topics

13 Projects

+100 000 contributors per month (english language)

How many page views will be seen by month ? 14 billion page views per month

When doing research on google, what is the chance to find a wikipedia article in the top 5 results ? In 96% case, top 5 results

# Wikipedia ranking :

# Wikipedia's quality labels:

| Class | Criteria | Reader's experience | Editing suggestions |
|---|---|---|---|
| ★ FA | The article has attained featured article status by passing an in-depth examination by impartial reviewers from WP:Featured article candidates.<br>**More detailed criteria** [show] | Professional, outstanding, and thorough; a definitive source for encyclopedic information. | No further content additions should be necessary unless new information becomes available; further improvements to the prose quality are often possible. |

| Importance | Criteria |
|---|---|
| Top | Core articles which are a "must have" for Wikiproject Software.<br>High-traffic articles which many people outside of computer software will already have a good knowledge of. |
| High | Most people interested in software will be familiar with the topic, and the article gives context to a number of other information software articles. Is mentioned by many books and academic papers, and discussed in detail in more than one. |
| Mid | Known to many people interested in software, even if not in detail. |
| Low | More specific and specialized content known only to some people interested in software.<br>Most individuals, standards and software projects will be low importance unless they are well known or have high adoption. |
| NA | Subject importance is not applicable. Generally applies to non-article pages such as redirects, categories, templates, etc. |
| ??? | Subject importance has not yet been assessed. |

| Stub | articles will fall into this category.<br>**More detailed criteria** [show] | insufficiently developed features of the topic and may not see how the features of the topic are significant. | Stub-class Article to step up to a Start-class Article is to add in referenced reasons of why the topic is significant. |

# Example :



Marks are obviously subjective and dependent of the contributor's culture. Also the rate differes from one project to the other.

From a reader point of view, it is difficult to understand the grade.

Some other ratings methods could be implemented such as
- Readers could evaluate each article (star rating)
- Using IA

# Why this article ?

- Machine learning methods to improve the precision of the quality measures of Wikipedia articles

- Innovative point of view : Using structure based information AND article based information

# Data set presentation :

**Article assessment method :**



**Data collection :**



Max(project's assessment)

Max(project's assessment)

| | |
|---|---|
| Number of *FA* articles | 2,415 |
| Number of *GA* articles | 3,160 |
| Number of *B* articles | 3,209 |
| Number of *C* articles | 3,322 |
| Number of *Start* articles | 4,110 |
| Number of *Stub* articles | 4,273 |
| Total | 20,489 |

TABLE III: Distribution of the data set within different quality classes

# Data set presentation :

**Data cleaning :**

❖ Remove two classes that were too small : categories A and B+

❖ Remove articles that have been deleted

**Data preprocessing :**

MEDIAWIKI

❖ Get Content

Wiki-class

❖ Compute structure-based features

TextStat

❖ Compute content-based

❖ Clean Data

# Feature selection :

*Hypothesis*: the writing style matters for measuring the articles quality.

| Structure-based features ⬛ Wiki-class | Content-based features ⬛ TextStat |
|---|---|
| ● Article length <br> ● Number of references <br> ● Number of outlinks to other Wikipedia pages <br> ● Number of citation templates <br> ● Number of non-citation templates <br> ● Number of categories linked in the text <br> ● Number of images / length of article <br> ● Information noise score <br> ● Article has an infobox or not <br> ● Number of level 2 headings/ Number of level 3+ headings | ● Flesch reading score (En) <br> ● Flesch-Kincaid grade level (US) <br> ● Smog index (En) <br> ● Coleman-Liau index (US) <br> ● Automated readability index (US) <br> ● Difficult words <br> ● Dale-Chall score <br> ● Linsear write formula (US military) <br> ● Gunning-Fog index |

# Content-based features (Examples)

❖ *A wikipedia infobox*



❖ $$flesch\_reading\_ease = 206.835$$
$$- (1.015 \times avg\_sentence\_len)$$
$$- (84.6 \times avg\_syllables\_per\_word)$$
$$(1)$$

❖ $$flesch\_kincaid\_grade = 11.8 \times avg\_syllables\_per\_word$$
$$+ 0.39 \times avg\_sentence\_len - 15.59$$
$$(2)$$

# Feature selection :

**Hypothesis**: the writing style matters for measuring the articles quality.

| Structure-based features   Wiki-class | Content-based features   TextStat |
|---|---|
| <ul><li>Article length</li><li>Number of references</li><li>Number of outlinks to other Wikipedia pages</li><li>Number of citation templates</li><li>Number of non-citation templates</li><li>Number of categories linked in the text</li><li>Number of images / length of article</li><li>Information noise score</li><li>Article has an infobox or not</li><li>Number of level 2 headings/ Number of level 3+ headings</li></ul> | <ul><li>Flesch reading score</li><li>Flesch-Kincaid grade level</li><li>Smog index</li><li>Coleman-Liau index</li><li>Automated readability index</li><li>Difficult words</li><li>Dale-Chall score</li><li>Linsear write formula</li><li>Gunning-Fog index</li></ul> |

**Question**: Several readability scores related ?

# Solutions proposed in the article :

| Algorithms | Hyper-parameters | Specificities | Accuracy |
|---|---|---|---|
| **Linear regression** | None | - dependent variable: quality class<br>- independent variables: the features<br> - converted the quality class to an integer: Stub to 0, Start to 1, C to 2, B to 3, GA to 4 and FA to 6 | 25% |
| **Multinomial logistic regression** | None | Standard | 60%<br>(5-fold Cross-Validation) |
| **KNN** | K = 3 | Using the Euclidean distance | 55% (5-fold Cross-Validation) |
| **CART** | None | Standard | 48% |
| **SVM** | None | Standard | 61% (5-fold Cross-Validation) |
| **Random Forest** | None | Applied uniquely on the structure-based features | 58% (5-fold Cross-Validation) |
| **Random Forest** | None | Applied on the complete set | 64% (5-fold Cross-Validation) |

# The most accurate model: Random Forest

**Performances:**

- 3 metrics:

    - Accuracy: 64%

    - AUC (Area Under Curve): 0,91

    - NDCG score: 0,987

# Pros & cons of these solutions :

## Pros

- Over-fitting problem taken avoided by the 5-fold cross validation

- The model improved the accuracy of Wikipedia quality prediction

- This paper provides advices for authors to improve the quality of Wikipedia articles

## Cons

- Only for english Wikipedia articles

- The data used to evaluate the model could be reconsidered. Taking the maximum of each articles' marks is not the best practise.

- The data are manually labelled, so the marks are subjectives

# Data set :



**What can we say about learning on subjective evaluations :**

A paper is discussing this problem : _The Success and Failure of Quality Improvement Projects in Peer Production Communities_. It aims to evaluate several rating groups that contributed to wikipedia assessments.

Mains conclusions are :

- Some articles were not correctly assessed, contributor failed to apply the assessment criteria.
- Assessments made by WEP have been given by students and it results in lower rating/writing experienced. Some groups are more efficient to produce high rated articles.

# References :

- A [paper](#) named : Measuring Quality of Collaboratively Edited Documents: the case of Wikipedia done

- A [paper](#) named : The Success and Failure of Quality Improvement Projects in Peer Production Communities

- A tool with rank suggestions : [Wikirank.live](#)

- [Link](#) to English Wikipedia Quality Asssessment Dataset. Data used to perfom classification

- [Link](#) to wikipedia assessment method

Thank you for your attention

# Annexe

# Model Evaluation :

- Application of different classification methods with 5-fold cross-validation techniques

  - Dataset divided into five equal parts (5- fold).

  - Four parts used as a training set / remaining part as the testing set.

  - Process repeated five times, each part being used as a testing set alternately.

- A good practical technique for bias-variance trade-off in evaluating machine learning algorithms.