

Rerunning OCR: A Machine Learning Approach to Quality Assessment and Enhancement Prediction

Pit Schneider¹ and Yves Maurer¹

¹National Library of Luxembourg, Luxembourg

Corresponding author: Pit Schneider, pit.schneider@bnl.etat.lu

Abstract

Iterating with new and improved OCR solutions enforces decision making when it comes to targeting the right candidates for reprocessing. This especially applies when the underlying data collection is of considerable size and rather diverse in terms of fonts, languages, periods of publication and consequently OCR quality. This article captures the efforts of the National Library of Luxembourg to support those targeting decisions. They are crucial in order to guarantee low computational overhead and reduced quality degradation risks, combined with a more quantifiable OCR improvement. In particular, this work explains the methodology of the library with respect to text block level quality assessment. Through extension of this technique, a regression model, that is able to take into account the enhancement potential of a new OCR engine, is also presented. They both mark promising approaches, especially for cultural institutions dealing with historical data of lower quality.

Keywords

optical character recognition; quality assessment; enhancement prediction; candidate selection; machine learning; historical data; cultural institutions

I CONTEXT

In the context of its digitization program, the National Library of Luxembourg (BnL) started a first initiative in the optical character recognition (OCR) space in 2006. Back then the external scanning suppliers were in charge of performing OCR on the scanned historical newspapers, using various software solutions over the years. Although OCR is considered a largely solved problem for modern documents (Doermann and Tombre [2014]), it remains a non-trivial task for historical data. That is why the library always considered the resulting output to feature a quality standard that could be improved in the future, with means of continuing software advancements.

A BnL pilot project conducted by Maurer [2017] proposed a framework to rerun OCR using a contemporary engine, such as Tesseract (Kay [2007]). The method leverages a metric that compares the new and original output on the ratio of number of characters belonging to words found in a dictionary. Altogether, the related article described promising results, served as a proof of concept and marked the starting point for subsequent OCR initiatives.

Fast forwarding to the year 2020, a new project is initiated, aiming to build a new in-house OCR model, referred to as NEWMODEL in the rest of this article. It was trained on BnL data

and represents an improvement on the original OCR quality. A prerequisite for the application of NEWMODEL, however, is a method that is able to first assess the original OCR quality, without relying on any ground truth counterparts. In terms of terminology, this technique is referred to as *automatic* OCR quality assessment. The motivation for employing such an approach and making it a prerequisite is threefold. It enables:

1. The reduction of computation time through selective targeting of reprocessing candidates.
2. The collection of statistical insights, estimating the improvement in OCR accuracy.
3. The lowering of the risk of a potential accuracy reduction for a subset of the data.

(1)

1.1 Related Work

Although the motivations in (1) seem compelling, tackling the problem of automatic quality assessment has seen relatively limited attention. More research has been devoted to the field of OCR post-correction techniques (e.g. Schaefer and Neudecker [2020]). Additionally, there are a number of studies (e.g. Strien et al. [2020] and Hill and Hengchen [2019]) that aim to assess the impact of sub-optimal OCR quality on natural language processing (NLP) tasks, without necessarily measuring the quality itself. The general ideas for automatic quality assessment are diverse and do not seem to yield a clear winner method in terms of potential. Altogether, they can be split into three classes, depending on the used data source.

First, image based methods have probably seen the widest amount of research, but also turn out to be the most computationally expensive. By neglecting the possibly already existing OCR output and by exclusively looking at the source material, they try to analyze features that could impact OCR engines. For instance, Blando et al. [1995] inspect specific typeface properties, while Lu et al. [2020] aim to discover physical image distortions. Other examples include quantifying image degradation (Peng et al. [2015]) and estimating the amount of blur (Kieu et al. [2016]). Finally, a recent work from Singh et al. [2018] uses surrogate models to learn document quality based on ground truth images.

Next, a common OCR engine produces output that extends beyond the raw text, which is the source for another class of approaches. In particular, Gupta et al. [2015] label erroneous text bounding boxes based on their spatial distribution and geometry, as an indicator for OCR quality. Similarly, Springmann et al. [2016] use engine confidence scores to compare the recognition quality of different models.

This leaves the processing of the output text only, which is the main source of inspiration for the present work. Here, the involvement of dictionaries has been looked at by researchers. For instance, Alex and Burns [2014] compare output words to the most similar entry in a dictionary. A very different approach comes from Cavnar and Trenkle [1994], showing that n-grams can be successfully used to categorize texts by comparing them to n-gram based classification profiles. Moreover, work from Kulp and April [2007] and Taghva et al. [2001] show the development of rule-based *garbage token* detection systems. Finally, an original contribution from Salah et al. [2015] uses a secondary (reference) OCR engine to cross-align output results using a Support Vector Regression technique.

1.2 Data

Subject to the application of NEWMODEL are the approximately 102,000 historical newspaper issues, dating from 1841 to 1954. The newspaper articles are mostly written in German (*de*), French (*fr*) and Luxembourgish (*lb*). Their typography is more or less evenly split between Antiqua and Fraktur typefaces, rendering the data rather diverse (Figure 1).

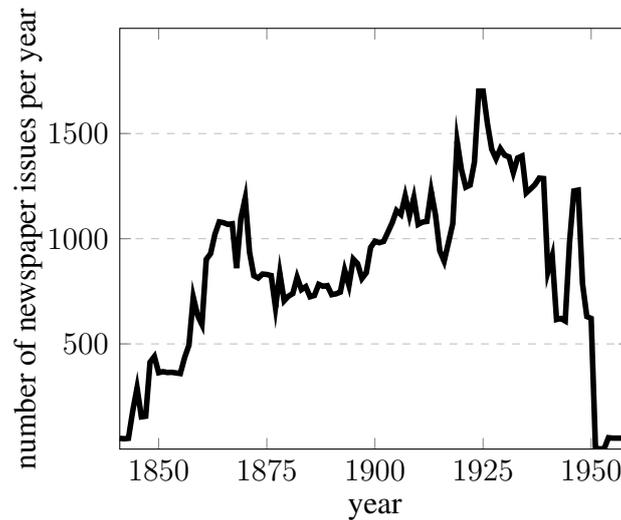


Figure 1: Publication date distribution of the BnL historical newspaper collection.

This work frequently refers to *blocks* as the most common data concept. A block represents the OCR output text, derived from the image of an individual paragraph or even a small article. As far as the layout is concerned, a block is always contained within a single text column. The choice to treat the data as a set of individual blocks is mainly motivated by the fact that there is a higher likelihood that properties, such as language or font, remain constant within a block.

The new OCR project has been initialized by first building a ground truth set. A subset of close to 7,000 text block images was selected and transcribed, mainly to serve for OCR training purposes. A second use case is to split-off of a separate test set, providing a foundation for automatic quality assessment. The possibility to test a given OCR output, by comparing it to its gold standard counterpart, is the basis for a supervised learning process. The resulting model will then be used for automatic quality assessment. Hence, finding a correlation between textual features, that can be computed without availability of a gold standard, and the text quality itself, is the venture discussed in the rest of this article.

From here on, this article will proceed by covering the two main contributions in the form of Sections II and III, before concluding by reflecting back on the statements in (1).

II QUALITY CLASSIFIER

This work proposes a machine learning based classifier that is designed to assess the text quality of an entire text block. With this intention, a couple of definitions need to be established first.

2.1 Definitions

Given the three motivations in (1), fitting a binary classifier with classes referring to *sufficient* and *insufficient* quality is the logical starting point. That is why the classes space C is defined as

$$C = \{0, 1\}, \quad (2)$$

with zero and one respectively referring to sufficient and insufficient quality. Coupling the positive class with bad OCR quality follows the notion of the classifier determining (a minority of) candidate blocks.

A supervised learning process is using training data T , given by

$$T = \{(X_1, Y_1 \in C), \dots, (X_n, Y_n \in C)\}. \quad (3)$$

It is also defined that every feature vector¹ has k dimensions, such that

$$X_i = (x_i^0, x_i^1, \dots, x_i^{k-1}). \quad (4)$$

The process of extracting all k features from i th text block B_i is referred to as the feature function

$$f : B_i \rightarrow X_i. \quad (5)$$

Having just defined B_i , there are a few additional notations related to a given block. While G_i is used to refer to the ground truth version of B_i , the cardinality $|B_i|$ returns the total number of characters (including whitespaces) within the block. Furthermore, B_i^t encodes all tokens (simple whitespace character delimitation) found in B_i . The concept of cardinality can again be utilized to obtain the length of a token. Lastly, the language function $\ell(B_i)$ returns the natural language of B_i .

The quality classifier can now be summarized as the function

$$\text{QUALITY} : B_i \rightarrow Y_i \in C. \quad (6)$$

To make QUALITY more robust it is trained on both the original OCR and NEWMODEL outputs, thus involving a variety of OCR software. This is illustrated in Figure 2, which also shows the high-level workflow and some of the just established notations.

¹Please note that in this article exponents are always used as labels and should not be interpreted as mathematical powers.

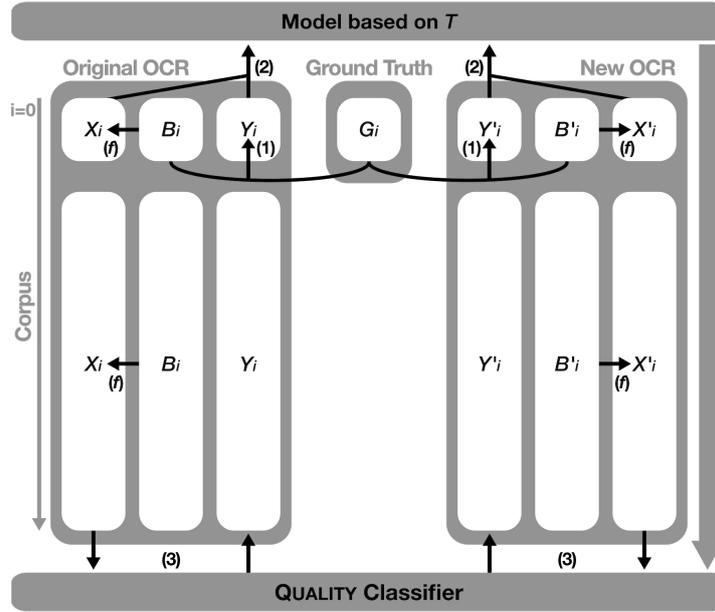


Figure 2: High-level workflow of (f) extracting text features, (1) determining the text quality, (2) training a model and (3) using that model to apply the classifier to the rest of the corpus.

2.2 Features

Next, focus is shifted to the topic of feature extraction. A collection of approaches and features has already been presented in Subsection 1.1. While sampling them for QUALITY, BnL considered features that focus on revealing shortcomings in the original OCR engine (e.g. Antiqua/Fraktur typeface confusion), and not necessarily in the source image quality. A second consideration aspect was computational efficiency, preventing significant impacts on the processing time of the new OCR pipeline. That is why the set retained for QUALITY is given by

- x^0 : dictionary mapping,
- x^1 : tri-gram comparison,
- x^2 : garbage token detection,
- x^3 : publication year consideration.

Thus, f operates on OCR output only. No feature is extracted from the source image, guaranteeing the efficiency of QUALITY.

2.2.1 Dictionary Mapping

A commonly used technique in automatic quality assessment is to compare the output words to a dictionary of the same language (e.g. Alex and Burns [2014]). Given a block B_i , its language $\ell(B_i)$, token $t \in B_i^t$ and dictionary $D^{\ell(B_i)}$, a binary variable is defined as

$$\begin{aligned} \text{map}(t, D^{\ell(B_i)}) &= 0 \text{ if } t \text{ is not in the dictionary,} \\ \text{map}(t, D^{\ell(B_i)}) &= 1 \text{ if } t \text{ is in the dictionary.} \end{aligned} \quad (7)$$

In the context of QUALITY, the feature x_i^0 is derived from B_i by computing ratio

$$x_i^0 = f(B_i)[0] = \frac{\sum_{t \in B_i^t} \text{map}(t, D^{\ell(B_i)}) \times |t|}{\sum_{t \in B_i^t} |t|}. \quad (8)$$

Given (8), every token is weighted by its own length, instead of simply returning the fraction of successfully matched tokens.

2.2.2 Tri-Gram Comparison

As suggested by Zipf [1949], given a natural language, the word rank-frequency distribution is a relation that is inversely proportional. The same law naturally holds for smaller entities within a language, such as n-grams. Building on this, Cavnar and Trenkle [1994] have demonstrated that texts can be categorized by comparing their contained n-grams to n-gram based classification profiles.

In a similar way, an n-gram similarity measure is established for QUALITY. More specifically, the measure makes use of the ranks of the top γ tri-grams in terms of frequency of language $\ell(B_i)$. The rank function $r(\text{tri}, \ell(B_i))$ returns the frequency rank of any tri-gram tri for language $\ell(B_i)$. Before computing the feature value, all possible character tri-grams are extracted from every $t \in B_i^t$. It should be noted that tri-grams are limited to only span across letter characters. For instance, there is

$$\begin{aligned} t \in B_i^t &= \text{Luxemb0urg} \\ \text{tri-grams for } t &: \{\text{lux}, \text{uxe}, \text{xem}, \text{emb}, \text{urg}\}. \end{aligned} \quad (9)$$

Let B_i^{tri} denote the set of all tri-grams in B_i . The feature value x_i^1 is calculated by

$$x_i^1 = f(B_i)[1] = 1 - \frac{\sum_{\text{tri} \in B_i^{\text{tri}}} \min(\gamma, r(\text{tri}, \ell(B_i)))}{\gamma \times |B_i^{\text{tri}}|}. \quad (10)$$

As a result of the exponential nature of the Zipfian distribution, the value of γ seems rather inconsequential as long as it is not too small. During the implementation process, $\gamma = 1000$ was chosen by BnL, safely covering all major tri-grams (in terms of importance) in the language. Naturally, the potential of this feature is increasing as $|B_i|$ increases as well.

2.2.3 Garbage Token Detection

As stated by Wudtke et al. [2011], a more serious category of OCR quality issues is the presence of tokens which also render the prediction of their correct replacement tokens infeasible. A feature, describing the amount of *garbage tokens* within B_i , combines ideas by Kulp and April [2007] and Taghva et al. [2001] into a set of nine precise rules.

A token $t \in B_i^t$ is identified as garbage in case it holds that t contains at least one of the following:

1. twenty-one characters.

2. three consecutive occurrences of the same character.
3. four consecutive vowels.
4. six consecutive consonants.
5. one vowel and at least one consonant and the count of one of them is more than eight times greater than the other.
6. one lower-case letter and even more upper-case letters.
7. one upper-case letter and starts and ends with a lower-case letter.
8. one alphanumerical character and contains even more non-alphanumerical characters.
9. two distinct non-alphanumerical characters, excluding the first and last character.

Applying the logical *OR* operator to this enumeration, a binary variable for token t is given by

$$\begin{aligned} \textit{garbage}(t) &= 0 \text{ if no rule applies,} \\ \textit{garbage}(t) &= 1 \text{ if at least one rule applies.} \end{aligned} \quad (11)$$

Hence, feature x_i^2 is extracted from B_i using

$$x_i^2 = f(B_i)[2] = 1 - \frac{1}{|B_i^t|} \sum_{t \in B_i^t} \textit{garbage}(t). \quad (12)$$

2.2.4 Publication Year Consideration

Through BnL data analysis it emerged that the original OCR quality is to some extent sensitive to the period of publication. This property mainly exists due to changes in the OCR engine used and in the source document quality. A yearly basis has been chosen to discretize time. This seems to be the smallest possible time unit that effectively correlates to changes in OCR quality. Thus, there is

$$x_i^3 = f(B_i)[3] = \textit{year}(B_i). \quad (13)$$

2.2.5 Language Detection and Independence

It remains to be addressed that the usefulness of two out of the four features, namely x_0 and x_1 , depends on correct language detection. Unfortunately, multilingual data coupled with variable OCR quality renders this task very challenging. BnL tries to overcome this issue by operating on a smaller (text block) level, rather than processing entire articles or pages (with a higher likelihood of language changes). The *langid* software (Lui and Baldwin [2012]) is used as a fallback after having run B_i against a selection of stop words for lb . Also, blocks where $\ell(B_i) \notin \{de, fr, lb\}$ are discarded (no prediction), opting for a trade-off with a higher accuracy but a little less volume.

Although garbage token detection (x_3) is not sensitive to the correct detection, it definitely needs to be robust enough to be language independent within $\{de, fr, lb\}$. Since Kulp and April [2007] and Taghva et al. [2001] based themselves on English data, it was necessary to review the nine listed rules for the languages in question. After careful consideration, it seemed sufficient to reassess rule 3 and 4 by verifying that four consecutive vowels and six consecutive consonants are indeed very rare appearances in all three languages.

2.2.6 Feature Experimentation

Features that did not contribute to the classifier performance, but were tested by the library, are listed below:

- A feature stating the font class (Antiqua/Fraktur), derived from the source image.
- A metric encoding the value $|B_i|$. Testing was backed by the hypothesis that smaller blocks (mostly headlines) would generally have a lower x^0 value induced by the presence of a higher ratio of named tokens not found in $D^{\ell(B_i)}$.
- A property indicating $\ell(B_i)$ through one-hot-encoding for a predefined set of language classes.

2.3 Class Definition

Before a classification model can be created, every B_i needs to be assigned a quality class $Y_i \in C$ in T . Here, the popular Levenshtein *edit* distance (Levenshtein [1965]) is used to compute quality measure

$$q(B_i) = 1 - \frac{\min(|B_i|, \text{edit}(B_i, G_i))}{|B_i|}. \quad (14)$$

Applying threshold θ leads to the class definition of

$$\begin{aligned} \text{if } q(B_i) \geq \theta : & \quad Y_i = 1 \in C, \\ \text{else : } & \quad Y_i = 0 \in C. \end{aligned} \quad (15)$$

2.4 Implementation

After having established the computation of T , QUALITY can be fit using a machine learning algorithm. Two non-linear methods were selected for this purpose:

- K-nearest-neighbour (KNN), motivated by its rather easy implementation, as a first choice.
- A feedforward neural network (NN) in view of potentially training a more complex model with more flexibility in terms of hyperparameters.

The NN architecture showing the best results features two identical *relu* activated hidden layers with 16 nodes, each followed by dropout of 0.5. Output layer classification is done through *softmax*. Other hyperparameters include a learning rate of 10^{-4} and a batch size of 1.

2.4.1 Preprocessing

Data standardization is applied in the NN case, for every d from 1 to k , in a way that

$$x^d = \frac{x^d - \bar{x}^d}{\sigma}, \quad (16)$$

with \bar{x}^d representing the mean and σ the standard deviation. For KNN to guarantee equal importance among features when computing the distance vectors, the feature value ranges need to be equal. That is why better results are obtained through min-max normalization, i.e.

$$x^d = \frac{x^d - \min(x^d)}{\max(x^d) - \min(x^d)}. \quad (17)$$

2.4.2 Training and Testing

QUALITY tries to mostly tackle, although influenced by threshold θ , an imbalanced classification problem, with the negative class outnumbering the positive one. This not only makes evaluation of the classifier less trivial, but creates challenges to train on enough positive data points.

To perform data augmentation and to specifically combat the lack of positive examples, two NEWMODEL outputs are generated for every block in the ground truth set.

1. A *new best-effort* version, with NEWMODEL being regularly applied, is included in blocks set B_{new} .
2. A *bad* version, with NEWMODEL purposefully applying a model trained on a different font (generating worse results), is included in blocks set B_{bad} .

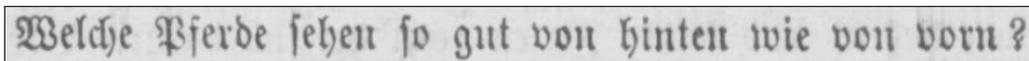


Figure 3: Very small example block source image.

Providing examples tied to Figure 3:

- "Welche Pferde sehen so gut von hinten wie von vorn?" $\in B_{new}$,
- "Welche Pferde sehen so gnt von hinten wie von vorn?" $\in B_{ori}$,
- "Belche serde fehen so gut von hinten wie von vorn?" $\in B_{bad}$.

The sets B_{new} and B_{bad} , together with the original OCR output B_{ori} , are contained within

$$B_{all} = \bigcup \{B_{new}, B_{ori}, B_{bad}\}. \quad (18)$$

Constant α (directly influenced by the chosen θ) is used to reference the positivity rate, which helps to quantify the imbalance of the problem. To provide an example, α_{ori} denotes the fraction of positive data points within B_{ori} . The set B_{all} forms the basis for a training/test set split. A fixed β blocks sized test set is first sampled from B_{all} by retaining positivity rate α_{ori} , thus creating a realistically imbalanced test scenario. The remaining blocks in B_{all} form the largest possible training set with respect to a perfect $\alpha = 0.5$ rate. In the NN case, 20% of the training set is split-off for validation purposes.

To evaluate QUALITY, next to the F_1 score (harmonic mean of precision and recall), emphasis is put on Cohen's Kappa (Cohen [1960]) metric, which takes class imbalance into account by

returning

$$kappa = \frac{p_0 - p_e}{1 - p_e}. \quad (19)$$

In (19) p_0 encodes the accuracy of the test set and p_e is the agreement between the model predictions and the actual class values, as if happening by chance (random class assignment).

2.4.3 Results

Results in Figure 4 are based on $|B_{all}| = 20,166$ and $\beta = 1,000$. Changes in α_{ori} do not seem to affect performance of the classifier significantly, pointing to a rather successfully handled class imbalance. The results can be seen as encouraging, but certainly still leave room for improvement.

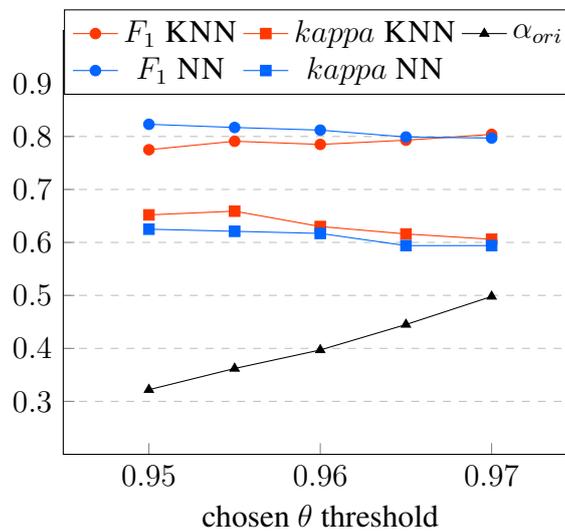


Figure 4: F_1 and $kappa$ scores of QUALITY given α_{ori} and θ .

A hypothesis coupled to experiments conducted with QUALITY, which potentially explains part of the model errors, states:

The quality class of smaller sized blocks (e.g. Figure 3) is considerably harder to determine. This is driven by the reduced amount of data for feature extraction. (20)

Therefore, the next section's results will make use of this observation by applying a weighted metric.

III ENHANCEMENT PREDICTION

While QUALITY incorporates a promising start to target an OCR rerun, it does involve a fundamental problem. More specifically for the BnL use case the downside of QUALITY lies in the lack of enhancement prediction, considering NEWMODEL. Classifying a block as insufficient does not imply that reprocessing changes the class, or even improves the quality at all. Moreover, a binary classifier is prone to provide limited feedback in terms of quality improvement

insights. Class conversions alone are not sufficient to obtain a good estimate on the overall improvement of the data.

3.1 Regression Definition

Based on this observation, a regression model is leveraged to compute enhancement predictions based on X_i . An adequate model naturally needs to output an estimate expressed in the same unit of measure as q , as defined in (14), entailing that

$$\text{ENHANCE} : B_i \rightarrow [-1, 1]. \quad (21)$$

To implement the regression model, T requires one modification. While X_i is left untouched, Y_i is replaced by a continuous variable. Therefore let i and j , with $i \neq j$, denote indices in B_{all} . Based on this, all block pairs are enumerated such that i and j reference the same source image and it also holds that

$$B_i \in B_{ori} \text{ and } B_j \in B_{new}. \quad (22)$$

Using (14), Y_i is computed in a way that

$$Y_i = q(B_j) - q(B_i) \quad (23)$$

now encodes the *potential* of the application of NEWMODEL, representing an information that is more valuable to the library while envisioning an OCR rerun.

3.2 Results

The machine learning algorithm with the best result is a regression version of KNN, returning the weighted (based on $|B_i|$) mean of all K neighbours. Applied on T , KNN outperforms other implementations, such as the same NN architecture (Subsection 2.4) with adjusted output layer and activation functions, or linear and logistic regressions.

3.2.1 Performance

To evaluate ENHANCE, the mean average error (MAE) measure by Willmott and Matsuura [2005] is used. This provides the ability of interpreting the model performance in the unit of measure q . The selected testing method of leave-one-out cross-validation was motivated by the fact that $|T|=6723$ is relatively small, making it computationally feasible to obtain a robust test result.

Considering $K=43$ neighbours, $\text{MAE}=0.034$ is achieved. This can be interpreted as a promising result, given that the test set features a high variance, more precisely a standard deviation of 0.14. Another reassuring aspect is that the model is only slightly too optimistic, by predicting 0.0029 too high on average. Overall, no fundamental bias can be observed.

As stated in (20), predicting on smaller blocks seems to be harder. This hypothesis can be reinforced by evaluating on an adaptation of MAE (here denoted as MWAE), which weights the loss (absolute difference between actual/predicted enhancement) of B_i by $|B_i|$. Since the size of

the block obviously directly correlates with the amount of text that is enhanced (or degraded), one can argue that MWAE even represents a fairer evaluation of ENHANCE. After all, a clear regression performance improvement comes with $MWAE = 0.024$ for $K = 31$.

3.2.2 Analysis

Reprocessing candidate selection based on ENHANCE requires a cut-off value, here again denoted as θ . Using the policy that every B_i , where it holds that

$$\text{ENHANCE}(B_i) \geq \theta, \quad (24)$$

is selected as a candidate (non-candidate otherwise), three ratios with respect to the total number of blocks, remain of particular importance:

- The ratio of candidates featuring a strict reduction in q , denoted as ϵ_r .
- The ratio of non-candidates featuring a strict increase in q , denoted as ϵ_i .
- The ratio of candidates, denoted as c .

The three ratios (calculated using weighting based on $|B_i|$) are depicted in Figure 5 for select values, such that $-0.06 \leq \theta \leq 0.16$. The graph shows a strong accuracy of NEWMODEL itself (rather low and flat slope of ϵ_r) and ratio ϵ_i properly adjusting to changes in θ .

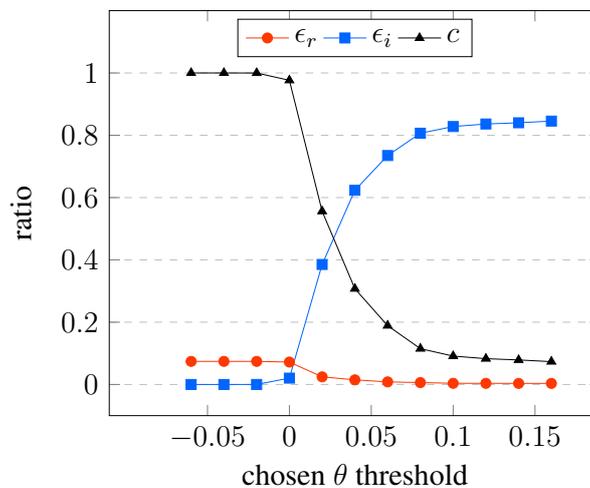


Figure 5: Ratios ϵ_r , ϵ_i and c for given values of θ .

Figure 5 shows that θ values, satisfying $0 \leq \theta \leq 0.05$, are most suitable for the application of ENHANCE, considering the BnL data and NEWMODEL.

A Python implementation and data model of ENHANCE, being part of the source code of the entire OCR project, can be publicly [accessed](#)².

²URL: <https://github.com/natliblux/nautilusocr>

IV CONCLUSION

This article commenced by enumerating three reasons to motivate the requirement of automatic quality assessment. To live up to those needs, feature extraction was discussed by looking at different ideas coming from the literature. Machine learning was applied to build the QUALITY classifier, designed for block level OCR candidate selection. Finally, this approach was extended through ENHANCE by considering the potential of NEWMODEL.

The three motivations in (1) can now be re-evaluated as follows:

1. At the time of writing, BnL already made use of QUALITY to save processing time. Using $\theta = 0.95$, leading to an appropriate balance in terms of target quality and reprocessing volume, a first experiment was conducted. Important to note is the 566,000 newspaper pages (102,000 issues) were processed in merely 15 days. This was enabled by the processing time of QUALITY, which generally stays below 5% of the time needed for the application of NEWMODEL itself.
2. Next, without statistical insights, the OCR rerun is comparable to a black box. It seems rather unfortunate for new artificial intelligence projects to enable better access to historical data, if those initiatives can not be advertised with concrete numbers. A first BnL application of QUALITY showed a class change for 70% of the candidate text lines. Additionally, deep diving into the feature values revealed positive average increments for x^0 , x^1 and x^2 . However, since this seems insufficient for a very clear picture, ENHANCE has been developed, expressing its predictions in the most comprehensible unit of measure, being q .
3. Lastly, ENHANCE reinforces the reduction of risks. The selection of candidates for reprocessing based on QUALITY is exposed to the risk of a poorly performing NEWMODEL. This problem is solved by ENHANCE, which can be applied using any cut-off threshold, depending on the amount of desired risk.

Altogether, ENHANCE will represent a very helpful addition to the newly developed OCR pipeline of the library and will serve as the basis for future reprocessing candidate selection processes.

The work described in this article has shown that estimating text quality and its potential to improve is a rather difficult task in itself, especially when computational efficiency without source image processing is desired. This is joined by the hurdles of language recognition, the availability of dictionaries covering historical language changes and the challenges involving smaller blocks. Nevertheless, a concrete, applicable and working solution has been proposed. That is why this article was redacted with the intention to share those findings with other cultural institutions with similar requirements.

References

- B. Alex and J. Burns. Estimating and rating the quality of optically character recognised text. *ACM International Conference Proceeding Series*, 2014.
- L.R. Blando, J. Kanai, and T.A. Nartker. Prediction of ocr accuracy using simple image features. page 319, 1995.
- W.B. Carnar and J.M. Trenkle. N-gram-based text categorization. *Ann Arbor MI*, 1994.
- J. Cohen. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20:37–46, 1960.
- D. Doermann and K. Tombre. *Handbook of Document Image Processing and Recognition*. Springer Publishing Company, Incorporated, 2014.
- A. Gupta, R. Gutierrez-Osuna, M. Christy, B. Capitanu, L. Auvil, L. Grumbach, R. Furuta, and L. Mandell. Automatic assessment of ocr quality in historical documents. page 1735–1741, 2015.
- M. Hill and S. Hengchen. Quantifying the impact of dirty ocr on historical text analysis: Eighteenth century collections online as a case study. *Digit. Scholarsh. Humanit.*, 34:825–843, 2019.
- A. Kay. Tesseract: an open-source optical character recognition engine. *Linux Journal*, 2007.
- V. Kieu, F. Cloppet, and N. Vincent. Ocr accuracy prediction method based on blur estimation. pages 317–322, 2016. doi: 10.1109/DAS.2016.50.
- S. Kulp and K. April. On retrieving legal files: Shortening documents and weeding out garbage. Special Publication 500-274, 2007.
- V. Levenshtein. Binary codes capable of correcting spurious insertions and deletions of ones. *Problems of Information Transmission*, 1:8–17, 1965.
- T. Lu, D. Ilic, and A. Doms. Noise characterization for historical documents with physical distortions. 11353: 77–87, 2020. doi: 10.1117/12.2559694.
- M. Lui and T. Baldwin. Langid.py: An off-the-shelf language identification tool. pages 25–30, 2012.
- Y. Maurer. Improving the quality of the text, a pilot project to assess and correct the ocr in a multilingual environment. *Relying on News Media. Long Term Preservation and Perspectives for Our Collective Memory*, 2017.
- X. Peng, H. Cao, and P. Natarajan. Document image ocr accuracy prediction via latent dirichlet allocation. pages 771–775, 2015.
- A.B. Salah, J.P. Moreux, N. Ragot, and T. Paquet. Ocr performance prediction using cross-ocr alignment. pages 556–560, 2015. doi: 10.1109/ICDAR.2015.7333823.
- R. Schaefer and C. Neudecker. A two-step approach for automatic OCR post-correction. pages 52–57, 2020.
- P. Singh, E. Vats, and A. Hast. Learning surrogate models of document image quality metrics for automated document image processing. pages 67–72, 2018. doi: 10.1109/DAS.2018.14.
- U. Springmann, F. Fink, and K.U. Schulz. Automatic quality evaluation and (semi-) automatic improvement of ocr models for historical printings. *arXiv: Digital Libraries*, 2016.
- D. Strien, K. Beelen, M. Coll Ardanuy, K. Hosseini, B. McGillivray, and G. Colavizza. Assessing the impact of ocr quality on downstream nlp tasks. 2020.
- K. Taghva, T. Nartker, A. Condit, and J. Borsack. Automatic removal of garbage strings in ocr text: An implementation. *The 5th World Multi-Conference on Systemics, Cybernetics and Informatics*, 2001.
- C.J. Willmott and K. Matsuura. Advantages of the mean absolute error (mae) over the root mean square error (rmse) in assessing average model performance. *Climate Research*, 30(1):79–82, 2005.
- R. Wudtke, C. Ringlstetter, and K. Schulz. Recognizing garbage in ocr output on historical documents. *Proceedings of the 2011 Joint Workshop on Multilingual OCR and Analytics for Noisy Unstructured Text Data*, 2011. doi: 10.1145/2034617.2034626.
- G.K. Zipf. *Human Behaviour and the Principle of Least Effort*. Addison-Wesley, 1949.