

Technic Report- Generating functions/trees for evaluating optimal binarizations

1st Crina-Andreea Gherase

Project Manager

*National University of Science and
Technology Politehnica Bucharest*

2nd Nicolae-Cosmin Cornea

Team Leader

*National University of Science and
Technology Politehnica Bucharest*

3rd Stefan Costin Manta

Software Developer

*National University of Science and
Technology Politehnica Bucharest*

4th Mara-Ilinca Mocanu

Software Developer

*National University of Science and
Technology Politehnica Bucharest*

5th Paul-Iulian Andrei

Tester

*National University of Science and
Technology Politehnica Bucharest*

Abstract—This report outlines a project that focuses on improving document digitization through enhanced binarization. It introduces two solutions, one for global and another for local binarization, aiming to intelligently combine algorithm outputs.

Index Terms—Global and Local Binarization, Image Processing, F-measure, Document Digitization.

I. INTRODUCTION

Historical documents suffer from various degradations due to ageing, extended use, several attempts of acquisition and environmental conditions [1]– [4]. The main artefacts encountered in historical documents are shadows, non-uniform illumination, smear, strain, bleed-through and faint characters. Those artefacts are problematic for document image analysis methods which assume smooth background and uniform quality of writing [1]. In handwritten documents, the writer may use different amount of ink and pressure and generate characters of different intensity or thickness, as well as faint characters. The same writer may write in different ways even within the same document. Similar problems, such as faint characters and non-uniform appearance of characters of the same font, are also encountered in historical machine-printed documents[5].

In image processing and document analysis, global binarization is a basic procedure that turns a grayscale image into a binary image in which pixels are categorised as background or foreground depending on a global threshold value. Global binarization aims to simplify the image and improve its contrast, which will facilitate the extraction of significant information for further analysis.

Global binarization is essential for addressing a number of issues related to document degradation in the context of historical document research. As was previously indicated, artefacts including shadows, uneven lighting, smearing, strain, bleed-through, and pale lettering are common in historical texts. Conventional document analysis techniques may find it

challenging to read the content effectively in light of these artefacts.

By creating a threshold that distinguishes the foreground (text or pertinent information) from the background in a way that is applicable to the entire document, global binarization approaches seek to address these issues. Choosing a suitable threshold is an essential stage since it affects the precision of the processing stages that come after. Various approaches, including Otsu’s method or adaptive thresholding, can be used, depending on the type of degradation and the document’s properties.

Global binarization techniques must be strong enough to manage such complications when dealing with historical handwritten texts, where variances in ink strength, pressure, and writing style are frequent. The goal is to produce a binary representation that faithfully represents the document’s content, enabling more accurate and efficient analysis for tasks like optical character recognition (OCR) or text extraction.

Overall, global binarization is a key preprocessing step in the analysis of historical documents, serving as the foundation for subsequent tasks aimed at restoring, preserving, and understanding valuable historical materials. Developing effective global binarization methods is essential for ensuring the success of document image analysis techniques applied to diverse and degraded document collections.

In contrast to global binarization, local binarization, also known as adaptive thresholding, recognizes that the optimal threshold value can vary across different regions of an image. This technique divides the image into smaller, localized regions and applies a threshold individually to each of these regions. The adaptive nature of local binarization allows it to better handle variations in intensity, background, and degradation within different parts of a document.

In the context of historical document analysis, local binarization becomes particularly valuable when dealing with documents that exhibit uneven illumination, faded ink,

or variations in writing characteristics. Traditional global thresholding methods may struggle to adapt to these local variations, leading to inaccurate segmentation and subsequent analysis.

Local binarization algorithms, such as Niblack, Sauvola, or Nick, calculate the threshold for each local region based on the statistical properties of the pixels within that region. These methods are effective in preserving fine details and ensuring accurate segmentation, especially when applied to documents with complex backgrounds or intricate writing styles.

The advantage of local binarization lies in its ability to address the challenges posed by diverse document conditions within a single image. By adjusting the threshold dynamically based on local characteristics, it helps maintain the integrity of the document's content and enhances the accuracy of subsequent analysis tasks, such as OCR or feature extraction.

In document image processing, combining global and local binarization algorithms in a multi-stage approach is a popular strategy. By combining the best features of both approaches, this hybrid strategy optimises the binarization process for old documents with a variety of degradation patterns and heterogeneous content.

To sum up, local binarization is an effective preprocessing tool that can be used to adapt to a variety of specific localised issues that aged and damaged documents bring. It is an essential part of the pipeline for historical document analysis. When included into workflows for document image analysis, it produces more accurate and significant outcomes for historical material preservation, restoration, and interpretation.

II. PROPOSED SOLUTION

Our Java application has as the main technical solution the class "GenerateTrees" that contains all the logic for building trees. The GenerateTrees class encapsulates the creation and evaluation of mathematical trees on a given dataset. The class contains five methods, each representing a distinct mathematical expression or tree structure.

A. First Equation

Calculates various statistical measures on specific positions of the dataset. Involves mean, square mean, product, power, harmonic mean, and other operations. The final result is the maximum of two intermediate results.

B. Second Equation

Computes geometric mean, product, square mean, minimum, and mean on different subsets of the dataset. Involves operations like multiplication by 2, $\frac{a - b}{a + b}$, and harmonic mean. The final result is the power of two intermediate results.

C. Third Equation

Utilizes random values to select subsets of the dataset for mean, product, and maximum calculations. Involves operations such as adding half, dividing by 2, $\frac{a - b}{a + b}$, square, geometric mean, and multiplication by 2. The final result is the sum of two intermediate results.

D. Fourth Equation

Computes products, minimum, geometric mean, sum, multiplication by 2, maximum, square, and power operations on specific positions of the dataset. The final result is the result of $\frac{a - b}{a + b}$ function for two intermediate results.

E. Fifth Equation

Involves square mean, harmonic mean, $\frac{a - b}{a + b}$, half, multiplication by 2, power, and maximum operations on values at prime and divisible-by-3 positions. The final result is the product of three intermediate results. Additionally, a utility method isPrim checks if a given number is prime.

F. Conclusions and remaining code

Overall, this GenerateTrees class provides a flexible framework for defining and evaluating complex mathematical expressions on input datasets, allowing for the generation of diverse trees based on statistical operations.

The code utilizes a separate Operations class for various mathematical operations, and logging is implemented using SLF4J.

Regarding CSV files, we parse them using a BufferedReader and insert the data into a HashMap structure. The key in this structure is the image identifier and the corresponding value is a list containing thresholds indexed by the algorithm type used. This is done using the class "CsvReader".

The "Testing" class is used to check the optimization scores for the generated trees by processing the global training data and matching it with values from the CSV file. The final output is the calculated average score for the file or a message indicating that a score could not be calculated.

The "RunTrees" class is responsible for running multiple instances of the mathematical tree generation process in parallel, using a thread pool. It is designed to execute various tree-generation methods defined in the GenerateTrees class on different datasets concurrently. The class utilizes an inner class, RunTree, to represent a single task for calculating a mathematical tree and writing the result to a file.

III. ARCHITECTURE

For our application our team chose to use a layered architecture which consists of the following layers:

A. Presentation Layer

Main Class: Responsible for starting and initializing the application. Contains the main method to kick off the application.

B. Business Logic Layer

GenerateTrees Class: Contains the logic for generating mathematical trees based on input datasets. Represents the core business logic of the application. Operations Class: Contains various mathematical operations used in tree generation.

C. Data Access Layer

CsvReader Class: Responsible for reading data from CSV files. Acts as the interface between the application and external data sources.

D. Application Layer

RunTrees Class: Orchestrates the execution of tree-generation tasks in parallel. Manages the interaction between the GenerateTrees class, the datasets, and the output files. **Testing Class:** Performs specific tests and calculations using the generated trees and input datasets.

E. Interaction Flow

The Main class will initiate the application by creating instances of the necessary classes. After that, the RunTrees class will manage the parallel execution of tree-generation tasks and will interact with the GenerateTrees class and the datasets. The GenerateTrees class implements the core logic for generating mathematical trees and will utilize the Operations class for mathematical computations. The CsvReader class handles the reading of data from CSV file. The Operations class provides reusable mathematical operations. The Testing class performs specific tests and calculations using the generated trees and datasets.

F. Relationships

The Main class serves as the entry point and initializes the application. The RunTrees class orchestrates the execution of tree-generation tasks and interacts with the GenerateTrees class. The GenerateTrees class contains the core logic for tree generation and relies on the Operations class for mathematical computations. The CsvReader class handles data access by reading information from CSV files. The Operations class provides a set of reusable mathematical operations. The Testing class performs specific tests and calculations using the generated trees and datasets.

This architecture provides a clear separation of concerns and promotes maintainability and extensibility. It also allows for easier testing of individual components and encourages good coding practices. Depending on the specific requirements and complexity of the application, the architecture can be adjusted or expanded.

IV. INTERMEDIATE RESULTS

The score for the file FirstTree 1s: 36.05179033165557
The score for the file SecondTree 1s: 49.68040139444379
The score for the file ThirdTree 1s: 45.38752908802357
The score for the file FourthTree 1s: 48.291424357404956
The score for the file FifthTree is: 33.00829992819355
We chose to move forward with the second, third and fourth tree.

V. CONCLUSIONS

In conclusion, our project addresses the critical issue of improving document digitization, particularly in the context of historical documents with various degradations. We proposed two innovative solutions: global binarization and local binarization, each tailored to handle specific challenges inherent in historical document analysis.

REFERENCES

- [1] A. Antonacopoulos and D. Karatzas, "Semantics-based content extraction in typewritten historical documents", Proc. Int. Conf. Document Anal. Recognit., pp. 48-53, 2005.
- [2] B. Gatos, I. Pratikakis and S. J. Perantonis, "Adaptive degraded document image binarization", Pattern Recognit., vol. 39, no. 3, pp. 317-327, Mar. 2006.
- [3] A. Antonacopoulos, D. Karatzas, H. Krawczyk and B. Wiszniewski, "The lifecycle of a digital historical document: Structure and content", Proc. ACM Symp. Document Eng., pp. 147-154, 2004.
- [4] A. Antonacopoulos and D. Karatzas, "Document image analysis for world war II personal records", Proc. Int. Workshop Document Image Anal. Libraries, pp. 336-341, 2004.
- [5] K. Ntirogiannis, B. Gatos and I. Pratikakis, "Performance Evaluation Methodology for Historical Document Image Binarization," in IEEE Transactions on Image Processing, vol. 22, no. 2, pp. 595-609, Feb. 2013, doi: 10.1109/TIP.2012.2219550.