# Automatically Inferring the Document Class of a Scientific Article

**Antoine Gauquier**
antoine.gauquier@ens.fr

Pierre Senellart
pierre@senellart.com
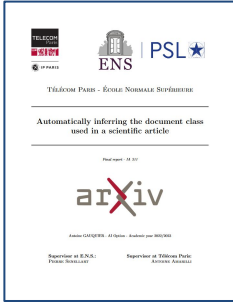
**23rd ACM Symposium on Document Engineering**

*Session 2 – Document Recognition, Summarisation and Inference*
*Wednesday, the 23rd of August 2023*

# Context

LaTeX source code

- Scientific style
- Mathematical formulas



Associated document in **PDF** format

Several libraries and commands.

Among them: `\documentclass{report}`

*(The argument defines the selected document class)*



Association for Computing Machinery (ACM)

`acmart`



American Astronomical Society (AAS)
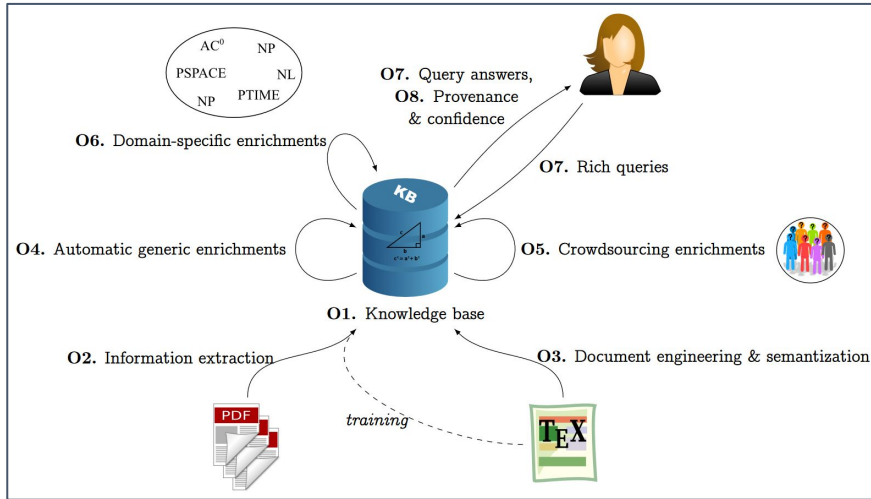
`aastex aastex61 aastex62`



American Mathematical Society (AMS)

`amsart amsproc`

1

# Applications

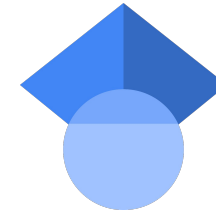## Systems extracting information from scholarly articles



The TheoremKB project
*https://github.com/PierreSenellart/theoremkb*

## Improve articles indexation in academic search engines



BASE search engine
*https://www.base-search.net/*



Google scholar search engine
*https://scholar.google.com/*

2

# Outline

**Dataset and performance metrics** ……………………………………………………………

**Statistical study** ……………………………………………………………………

**Random forest-based approach** …………………………………………………………

**CNN-based approach (deep learning)** …………………………………………………

# Dataset and performance metrics



Among these 98713 articles → more than 1200 document class names. We kept the most frequent ones, and merged the most similar ones (amsart and amsproc for instance), ending in 33 document classes.

$$\text{precision}_i = \frac{\text{TP}_i}{\text{TP}_i + \text{FP}_i}$$

$$\text{recall}_i = \frac{\text{TP}_i}{\text{TP}_i + \text{FN}_i}$$

$$\text{F}_1\text{-score}_i = 2\frac{\text{precision}_i \times \text{recall}_i}{\text{precision}_i + \text{recall}_i}$$

Why **Macroscopic F1-Score** ?

- Macroscopic gives same weight to each document class
- F1-Score gives finer analysis for multiclass classification than accuracy

4

# Statistical study
## Construction of five (simple) hand-designed features (1)

$$1 \text{ point} = \frac{1}{72} \text{ inch}$$

Weighted average left margin **(lm)**

$$\overline{m^{\mathrm{h}}} = \frac{\sum_{i=1}^{N_b} m_i^{\mathrm{h}} \times l_i}{\sum_{i=1}^{N_b} l_i}
\begin{cases}
m_i^{\mathrm{h}} & \text{Margin for i-th text block} \\
l_i & \text{Vertical height of i-th text block} \\
N_b & \text{Total number of blocks}
\end{cases}$$

Average first top margin **(tm)**

$$\overline{m^{\mathrm{v}}} = \frac{\sum_{i=1}^{N_p} \min_j m_{i,j}^{\mathrm{v}}}{N_p}
\begin{cases}
m_{i,j}^{\mathrm{v}} & \text{Distance between top of page of j-th block of i-th page} \\
N_p & \text{Total number of pages}
\end{cases}$$

5

# Statistical study

## *Construction of five (simple) hand-designed features (2)*

$$1 \text{ point} = \frac{1}{72} \text{ inch}$$

4 pts  2   136 pts   K. Papadopoulos and A. Syropoulos   6 pts

> Dynamical systems are characterized by equations that describe their evolution. A dynamical system is called *linear* when its evolution is a linear *process*. A process is linear when a change in any variable at some initial time produces a change in some variable at some later time, however, if the initial variable changes $n$ times, then the new variable will change $n$ times at the later time. In other words, any change propagates without any alterations. Any system that is not linear is called a *nonlinear* dynamical system [13]. A basic characteristic of these systems is that any change in a variable at some initial moment leads to a change to some variable at a later time, which is not proportional to the initial change. For example, the *logistic map* [12]

118 pts

334 pts

$$x_{n+1} = rx_n(1 - x_n),$$

91 pts   10 pts

> where $x_n \in [0,1]$ is the magnitude of population in generation $n$ and $x_{n+1}$ the magnitude of population at generation $n+1$, is a typical example of an equation that describes a nonlinear system. In this case, the system is the population of some species and the dynamics the changes from one generation

44 pts

334 pts

### Weighted average column width **(cw)**

$$\overline{w} = \frac{\sum_{i=1}^{N_b} w_i \times l_i}{\sum_{i=1}^{N_b} l_i}$$

$\begin{cases} w_i & \text{Width of i-th text block} \\ \\ l_i & \text{Vertical height of i-th text block} \end{cases}$

2   K. Papadopoulos and A. Syropoulos

> Dynamical systems are characterized by equations that describe their evolution. A dynamical system is called *linear* when its evolution is a linear *process*. A process is linear when a change in any variable at some initial time produces a change in some variable at some later time, however, if the initial variable changes $n$ times, then the new variable will change $n$ times at the later time. In other words, any change propagates without any alterations. Any system that is not linear is called a *nonlinear* dynamical system [13]. A basic characteristic of these systems is that any change in a variable at some initial moment leads to a change to some variable at a later time, which is not proportional to the initial change. For example, the *logistic map* [12]

$$x_{n+1} = rx_n(1 - x_n),$$

> where $x_n \in [0,1]$ is the magnitude of population in generation $n$ and $x_{n+1}$ the magnitude of population at generation $n+1$, is a typical example of an equation that describes a nonlinear system. In this case, the system is the population of some species and the dynamics the changes from one generation

### Most common font family **(ff)**
### Most common font size **(fs)**

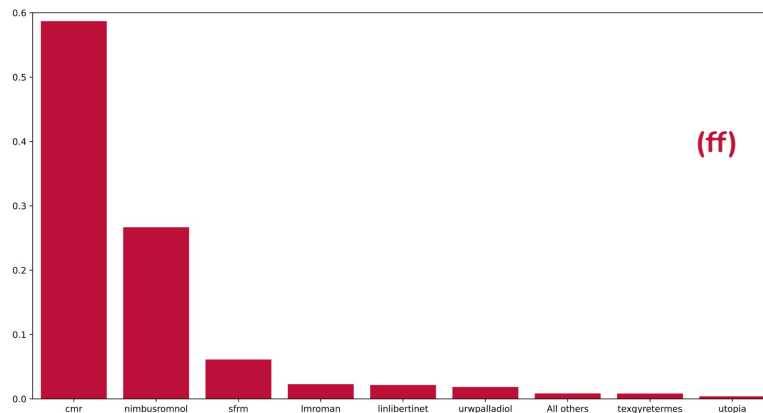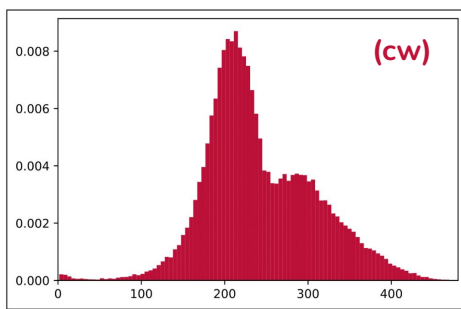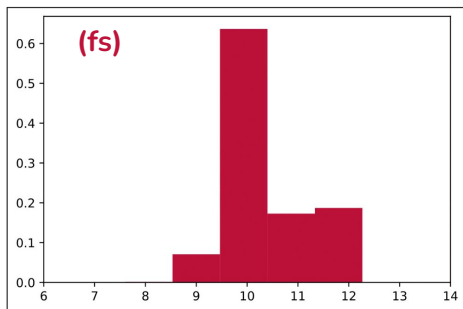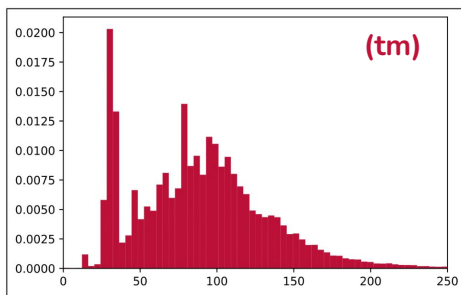$$f_i = \frac{\sum_{s \in S_i} l_s \times h_s}{\sum_{j=1}^{N_f} \sum_{s \in S_j} l_s \times h_s}$$

$\begin{cases} S_i & \text{Set of all tokens of i-th font} \\ \\ l_s & \text{Length of token s} \\ \\ h_s & \text{Height of token s} \end{cases}$

6

# Statistical study
*Global distributions of the features*



Only a few font families **(ff)** are widely used.

- Gaussian-like distributions for **(lm)**, **(tm)** and **(cw)**.

- Different values for **(fs)** and peak value for **(tm)**.

We identify several characteristics that seem to be document-class specific, and therefore **discriminative**.

# Statistical study

## Comparison of distributions from two different document classes



Source: https://doi.org/10.1051/0004-6361/201730392

bmvc2k document class

Source: https://doi.org/10.48550/arXiv.1801.09436

aa document class

This example shows that we can easily **separate** these two document classes with 3 features only.

**This entire study indicates that using statistical learning should work pretty well …**

# Random forest-based approach
## *Configuration and results*

Random forest model : Ensemble method that uses statistical learning to train a lot of decision trees on different subparts of the training dataset.

**Hyperparameters**
{
Minimum number of samples per leaf → set to 0.01% (risk of overfitting if too high)

Number of decision trees → set to 1000 to ensure stability of most common decisions
}

**Features of the model :** the five hand-designed features. **Output :** Predicted document class among 33 of them.

| Model | Averaged precision | Averaged recall | Macroscopic F1-Score |
|---|---|---|---|
| *Dummy* | *0.09%* | *3.03%* | *0.18 %* |
| Random forest model | 64 % | 66 % | 64 % |

Simple modelization (no deep learning and only five, simple, features) → Really promising results !

# CNN-based approach
## *Input data*

### Text element specific to AAS document class

DRAFT VERSION JANUARY 3, 2018
Typeset using LATEX **twocolumn** style in AASTeX61

A MODEL FOR DATA CITATION IN ASTRONOMICAL RESEARCH USING
DIGITAL OBJECT IDENTIFIERS (DOIS)

JENNY NOVACESCU,[1] JOSHUA E.G. PEEK,[1] SARAH WEISSMAN,[1] SCOTT W. FLEMING,[1] KAREN LEVAY,[1] AND
ELIZABETH FRASER[1]

*Source : https://arxiv.org/pdf/1801.00004.pdf*

### Example of input bitmap rendering



**256**

*Source : https://arxiv.org/pdf/1806.06252.pdf*

### Some usual elements from ACM document class



- ACM Reference Format
- Rights and information about the article

10

# CNN-based approach
## *Architecture*



Size: (256, 256, 1)

Size: (254, 254, 32)

Size: (63, 63, 32)

Size: (15, 15, 32)

Size: (3, 3, 32)

Size: 3 x 3 x 32

Flattening operation

**Convolutional layer**
Kernel size of (3, 3)
**ReLU activation**

**Convolutional layer**
Kernel size of (3, 3)
**ReLU activation**
**Max-pooling**
Kernel size of (4, 4)
**Dropout operation**
Probability 0.25

**Convolutional layer**
Kernel size of (3, 3)
**ReLU activation**
**Max-pooling**
Kernel size of (4, 4)
**Dropout operation**
Probability 0.25

**Convolutional layer**
Kernel size of (3, 3)
**ReLU activation**
**Max-pooling**
Kernel size of (4, 4)
**Dropout operation**
Probability 0.25

**Dense layer**
**Softmax activation**

0
1
2
32

# CNN-based approach
## *Results and comparison with state-of-the-art*

| Architecture | Macro F1-Score | Number of parameters | FLOPS (in billions) |
|---|---|---|---|
| *Our architecture* | *92.31 %* | ***38 177*** | *1.36* |
| ResNet50V2 | 92.28 % | 23 632 417 | 9.13 |
| NASNetMobile | 91.31 % | 4 304 597 | 1.50 |
| EfficientNetV2B0 | **93.43 %** | 4 091 844 | **0.80** |

## Analysis:

➔ 100 times less parameters than other models
➔ Almost as performant, above 92% of F1-Score
➔ Number of floating operations at inference time slightly above EfficientNetV2B0

# CNN-based approach
## *Separating heterogeneous document classes with reject option (1)*

| Document class | Precision | Recall | F1-Score |
|---|---|---|---|
| book | 56.84 % | 21.39 % | 31.09 % |
| report/wlscirep | 52.09 % | 77.69 % | 62.37 % |
| other (including article) | 69.17 % | 65.00 % | 67.02 % |

Common ground of theses classes : they are widely customizable, and thus embed a great **heterogeneity** of renderings.

What about directly putting apart these heterogeneous classes before applying classifier ? This is **reject option**.

# CNN-based approach
## *Separating heterogeneous document classes with reject option (2)*

| Model | Precision | Recall | F1-Score |
|---|---|---|---|
| Rejector | 90.55 % | 89.15 % | 89.04 % |
| Classifier | 96.94 % | 96.73 % | 96.83 % |

Improvement in the classifier performance (more than 4.50% in averaged F1-Score). However …

The rejector has lower performance → overall system not necessarily better !

Still very useful for applications where we know that **heterogeneous classes are not frequently observed or relevant** (for instance, articles from conference proceedings or journals).

*Recall for non-heterogeneous class of rejector is above 98 % : non-heterogeneous classes are almost always classified as so.*

# Conclusion and perspectives

- It is statistically relevant to discriminate document classes on the basis of features from PDF rendering.

- A (relatively) simple classification method on a set of 5 simple features gives promising results.

- Using a computer-vision based approach (CNN) gives really good performance, comparable to state-of-the-art models with way more parameters.

- We can even improve these results by putting apart heterogeneous classes, which are not related to a specific conference or journal.

- The experiment was conducted on a « small » subset of ArXiV (only 2018): what happens on a larger time frame?

- Dependency on ArXiV: we don't know any dataset where document class is readily available.

- We did show that using document class helps detecting mathematical environments (TheoremKB). But finding an efficient way maximising performance is still in progress.

# Thank you for your attention!
## Any questions?

https://github.com/AntoineGauquier/inferring_document_class_of_scientific_article/

**Antoine Gauquier**
antoine.gauquier@ens.fr

Pierre Senellart
pierre@senellart.com