
CONTENT CLASSIFICATION AND KEYWORD EXTRACTION FOR *ici.radio-canada.ca*

A PREPRINT

David Alfonso
RALI, Université de Montréal

Vincent Barnabé-Lortie
Médias Numériques, Radio-Canada

Shivendra Bhardwaj
RALI, Université de Montréal

Ilan Elbaz
RALI, Université de Montréal

Abbass Ghaddar
RALI, Université de Montréal

Fabrizio Gotti
RALI, Université de Montréal

Philippe Langlais
RALI, Université de Montréal

Guillaume Le Berre
RALI, Université de Montréal

Vincent Letard
RALI, Université de Montréal

Peng Lu
RALI, Université de Montréal

Laura Elisa Salas
DIRO, Université de Montréal

Olivier Salaün
RALI, Université de Montréal

Jason Jiechen Wu
RALI, Université de Montréal

December 17, 2019

1 Introduction

Radio Canada publishes between 450 and 600 articles per day on the *ici.radio-canada.ca* site. To this date, their database contains 901 156 articles. These are mandatorily annotated by the authors using 26 themes roughly corresponding to traditional paper columns and optionally annotated using 464 Sub-themes.

The themes and sub-themes help both the human reader and the search engine bots find information in the correct domain of interest, browse multiple news connected by themes and thematically limit the field of research of a specific article.

The journalists and content editors who are in charge of choosing the theme might not always agree on how to classify a particular article. Moreover, depending on the date, context and personal preference, one slightly ambiguous article can be classified in a non-evident Theme. Concerning the Subthemes, they are simply too many and finding the right one for each article becomes a tiresome task.

This is why *Radio-Canada* wishes to improve the consistency of their classification and reduce the amount of effort for journalists and content editors by proposing automatically detected Themes and Subthemes.

Radio-Canada also wished to enrich the content description of each article. Because, historically, the idea of having keywords was not introduced in the journalistic domain until very recently, in their current state, *ici.radio-canada.ca* articles do not include them.

For this reason we were asked to study the problem of automatically extracting keywords from the article content. The ideal situation would allow not only to offer keyword suggestions to journalists and content editors but to automatically label the 901 156 articles of the data base.

2 Pre-workshop data preparation

In order to avoid the same problems encountered on the previous workshop (where some time was lost in getting and distributing the industrial data) to collect and prepare data in collaboration with *Radio-Canada* (in the person of Vincent Barnabé-Lortie), Fabrizio Gotti heavily worked on it during the week previous to the workshop so the rest of the team could get to work quicker. He facilitated the acquisition of the data, saved it in a convenient format on a server of RALI, developed a dedicated Application Programming Interface (API) for efficiently access the data. He also acquired some potentially usefull data from journal *Le Devoir* website (around 12 000 of news with their respective keywords).

3 Data analysis

The first team task we needed to tackle was the analysis of the data. Before we could translate the industry problems into implementable tasks and sub-tasks, we needed to look at the data and understand what were the tools and raw material at our disposal.

After exploring the data a bit, Vincent Barnabé-Lortie, the *Radio-Canada*, *Médias Numériques* representative, lighten us on the doubts that arose concerning both the structure of the data / meta-data and the journalists' procedures when submitting an article. The main ones being:

- Each article mainly consists of a title, a summary, a lead paragraph, the body (all the paragraphs), the theme and, optionally, the sub-theme.
- Some articles may be empty for some of the sections mentioned above, this mainly happens in the case of news tickers, which do not have the traditional structure of an article, per se.
- Some of the themes and sub-themes in the taxonomy are deprecated and no longer used, even though they do appear in the data corpus.
- The journalists and content editors are required to specify one theme (and only one) for the article among the existing options.
- One of the themes in the taxonomy is called "Aucun thème sélectionné" ("No theme selected") and it seems to be used for old articles preceding the theme/sub-theme taxonomy that were not labelled back when they were written.
- It is not mandatory for the journalists and content editors to specify a sub-theme, but there is virtually no limit on the number of sub-themes an article may contain.
- There is no procedure that allows the journalist or content editor to propose a list of keywords for each article.
- One of the tools provided by *Radio-Canada* was a dictionary linking sub-themes to themes, but we observed that this theme/sub-theme link was not always respected (as further shown in Section 4).

4 Some statistics

Before continuing with the description of the tasks, we wish to rapidly expose some statistics on this specific data, which served us as guidance in the resolution of the exposed problems.

The observations derived from Figures 1, 2, 3, 4 and 5 lead us to formulate how we needed to clean the data.

As shown on Figure 1 some themes are not active or very little (and the sub-themes even more so). This means that the distribution is unequal and possibly biased towards certain themes more than others, as we will observe in the following Sections.

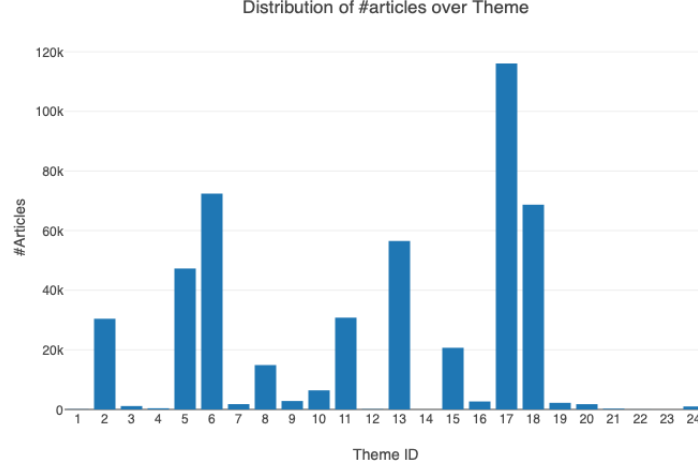


Figure 1: Distribution of the number of articles over the theme, showing some themes are not active in the whole data set.

On Figure 3 we observe that the sub-themes labeling is very imbalanced. We hypothesize that this tendency is a reflect on the need for popular sub-themes to distinguish themselves from the more general themes, which would create a motivation for the journalist to label the article with a non-mandatory sub-theme. We observe that the top 3 most frequent sub-themes are *Hockey*, *Politique provinciale* and *Éducation*.

It is worth noting that some sub-themes are not linked to one exclusive theme but to several themes (e.g.: the sub-theme *Mental health* is often linked to the themes *Health*, *Society*, *Miscellaneous news (Faits divers)*, etc.). This is further analyzed in Figures 4 and 5.

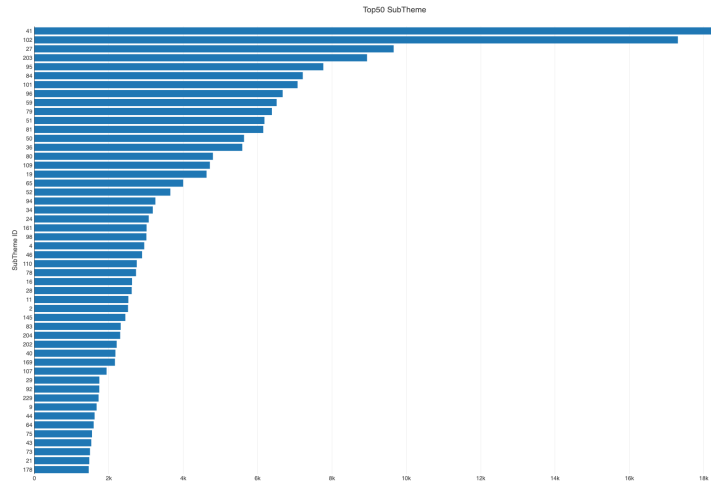


Figure 2: Excerpt of the distribution of the top 50 sub-themes over the number of articles, showing the most common sub-themes in the data set.

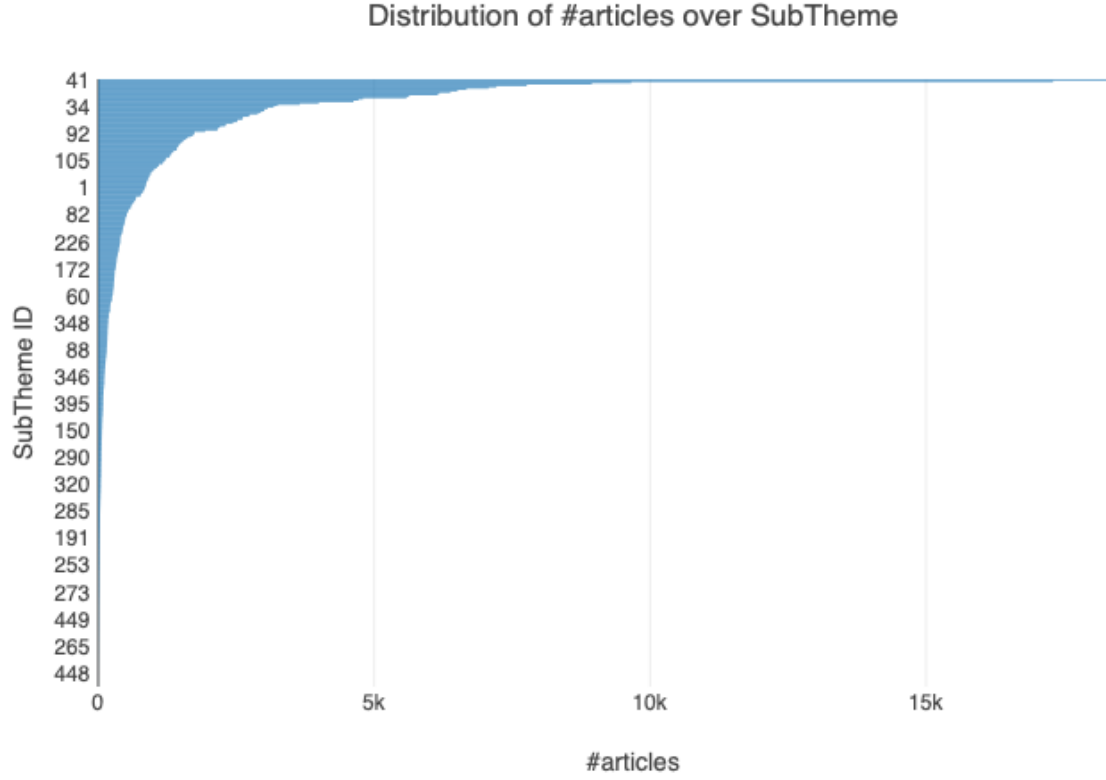


Figure 3: Sorted distribution of the sub-themes over the number of articles, showing that the distribution of sub-themes is also imbalanced.

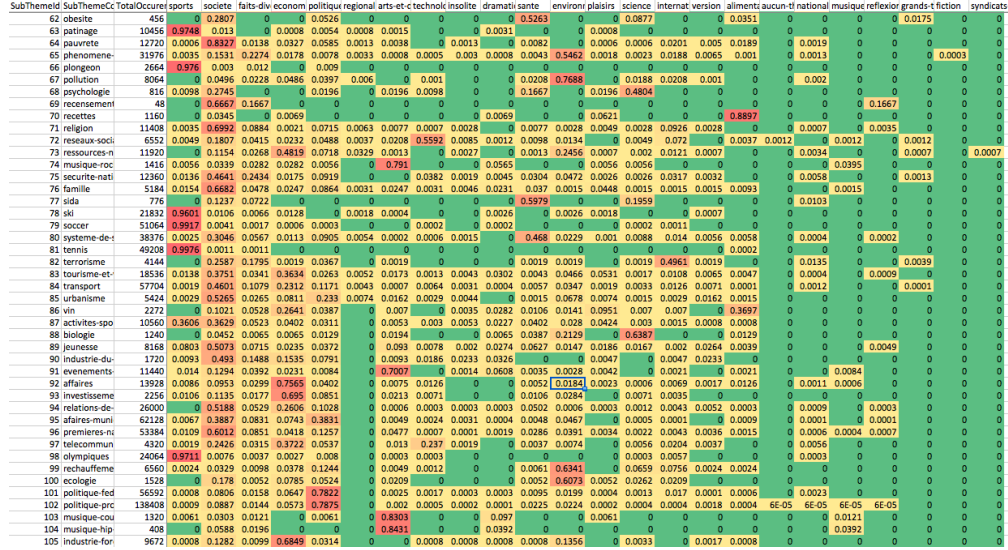


Figure 4: Excerpt of the heatmap of themes (horizontal) and sub-themes (vertical) correspondence, showing sub-themes spread over numerous themes. This indicates that some sub-themes are not limited to specific themes.

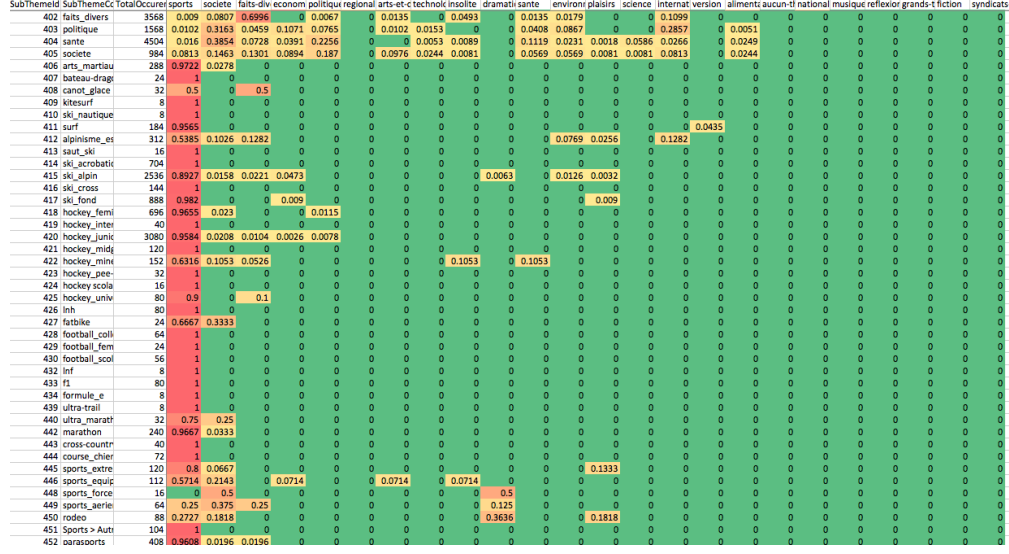


Figure 5: Excerpt of the heat-map of themes (horizontal) and sub-themes (vertical) correspondence, showing sub-themes very focused on one of very few themes. This indicates that some sub-themes are rather "faithful" to a specific theme.

5 Data cleaning and splitting

As a group, we agreed on the need to clean the data and use the same subsets to train our models and evaluate ourselves if we wanted our results to be comparable.

As we said on Section 3, the sub-theme is not a mandatory requisite for publishing an article on *ici.radio-canada.ca*. This means that a lot of the articles have a theme but no sub-theme. In order to work on equal ground when using themes and sub-themes classifiers, we decided to remove all articles that did not contain both a theme and (at least) one sub-theme.

There is one theme that corresponds to "No theme selected". Since this null-theme matches no specific domain, we chose to remove all articles labeled with this theme from our data set.

After cleaning, the total number of articles fell from around 900 000 to, approximately, 240 000. After this, the remaining data was split into 3 parts: the Training set, the Validation set and the Test set. After having randomized the whole data set, of the 240 000 total articles, 80% were assigned to the Train set, 10% to the Validation set and another 10% to the Test set. This was made once and distributed to the whole team so we could train, validate and test each model on the same subsets and obtain comparable results.

We also decided of a specific output format of the various classifiers we tested and an evaluation protocol for which an algorithm was implemented in order to assure evaluation uniformity for every team.

These tasks mainly fell on Fabrizio Gotti's shoulders, since they were necessary to start working on the problems as soon as possible and he already had an advance knowledge of the data, the format and the API.

Despite one week for preparing data before the workshop, the team spent more than a day to decide of a way to clean the data, and to proceed.

Table 1: Classification task scores.

Sub-task	Classifier	Precision @1	Recall @1	F1 @1
Theme identification	BERT	75.4%	75.4%	75.4%
Theme identification	Fasttext	77.6%	77.6%	77.6%
Theme identification	Logistic regression	77.8%	77.8%	77.8%
Theme identification	SVM	78.7%	78.7%	78.7%
Sub-theme identification	BERT	7.1%	5.8%	6.4%
Sub-theme identification	Fasttext	71.1%	58.1%	63.9%
Sub-theme identification	Logistic regression	na	na	na
Sub-theme identification	SVM	na	na	na

6 Theme/Sub-theme classifiers

The first problematic consisted in training in a supervised way some classifiers to predict the theme or the sub-themes on an article, based on its contents, that is, its words.

Due to time constraints, we tested two well-established feature-based approaches, namely Linear Regression and Support Vector Machines (Cortes et Al. [2]), as well as two approaches based on deep learning, namely BERT (Devlin et al. [3]) and Fasttext (Joulin et al. [4]). It is important to stress that many of those algorithms are controlled by hyper-parameters that we did not have time to investigate.

All resulting scores are shown on Table 1.

6.1 Logistic Regression

We first experimented with logistic regression, where each article is represented by a vector of dimension V , where V is the vocabulary size (the size of the list of each word in the training material). The value of each coefficient is computed with the standard TF-IDF score; very roughly frequent words in a document get rewarded by this score if not too frequent in the other articles. Training such a model amounts to learn to weight each dimension (word) of the associated vector space representation so to minimize error classification. The advantage of such a technique is that we can observe the learned weights and see what the model thinks about the importance of a specific word for a given theme (or sub-theme). We used Scikit-Learn for training the model nearly in its default setting. A modest number of variants have been tried that rendered very close performances in the end.

6.2 SVM

Support Vector Machine has been repeatedly reported as a robust classifier. Therefore we tested one such approach. The feature representation was very similar to the one given to the linear regression approach. We used the implementation in Scikit-Learn and its SGDClassifier function with almost all of the default arguments except for the random state value (42), the maximum of iterations (limited to 16), and tolerance (None). This model rendered the best performance as a theme classifier.

Both the Logistic Regression and SVM models can run on ordinary CPUs and do not require a state-of-the-art computation power.

6.3 Bert

BERT is a pretrained model. It provides contextual embeddings of the words in a sentence and can ultimately be used to extract a vector representation of a sequence of tokens. The model can be further fine-tuned to obtain better representation for a particular task.

For the theme/sub-themes classification task, we fed to BERT the first 128 tokens of the body and we took the embedding of the first token (corresponding to the "Start" symbol) to be our sentence representation. This representation

is then fed to two distinct linear layers followed by softmax layers (for themes and sub-themes). We used a cross-entropy loss for training. For sub-themes, experiments have been made using a multi-label objective but this method showed really weak results and we finally changed it to be a cross-entropy loss using only the first sub-theme as gold label.

Training is done using Adam with a learning rate of 10^{-5} . For these experiments we used a pytorch implementation of BERT (pytorch-transformers) and the "bert-base-multilingual-cased" model.

6.4 Fasttext

Another popular and recent neural network pre-trained model we tested was Fasttext, from Facebook's AI Research lab.

We used the basic available model and chose, as hyper-parameters, a learning rate of 1.0, using a vector dimension of 50, a word ngram length of 2, a bucket value of 200 000 and we trained for 40 epoch. As a result, we got a good score for theme classification and the best score for sub-theme classification. It is noteworthy to say that these sub-theme results rendered the best scores for the sub-theme classification sub-task.

Contrary to the Logistic Regression and SVM models, the BERT and Fasttext models require to run on highly potent GPUs in order to get a relatively fast result. BERT was trained using four Nvidia RTX 2080 Ti and convergence is achieved after approximately 2 hours. BERT was trained using one Titan XP and convergence is achieved after approximately 4 hours.

6.5 Future work

One of the most evident conclusions the classifier task teams arrived to, is that in order to better understand the data, we need to make a better analysis of the labeling done by the journalists. Understanding where are their priorities when they label an article with a specific theme or sub-theme is key to understand how clean or uniform the data is.

For the theme classification task, the best performing model was SVM. We expected it to outperform the Linear Regression model, but we were more surprised by the lower performance of the deep learning approaches. Several reasons can explain this, among which the relatively small training set, and the multilingual embeddings we used to seed the model.

As we can observe, there is room for improvement yet. Even with our best efforts, there is only so much we can do in just one week of work with the data. Our highest F1 score for theme classification is 78.7% and 63.9% for sub-theme classification. We believe these scores can be improved by testing other models, better adjustment on the hyper-parameters or even an automatic data cleaning for ambiguously labeled themes.

7 Keyword extraction

In order to allow an improved article access to the potential reader, the theme is not always enough. This is because themes and sub-themes only represent a more or less general thematic domain while, very often, the reader has a greater interest in a specific subject. This is where keywords become useful.

Since the database does not have any keyword annotations and we lack a large number of reader-annotators to manually extract, pinpoint or label each one of the 900 000 articles with multiple keywords, we moved to an automatic keyword extraction method. By which we mean to analyze the content of the article in order to deduce what words (or group of words) would be representative keywords.

The greater problem we encountered with this task, was the evaluation side of it. Without a keyword annotated corpus, we were unable to produce any evaluative test other than the superficial human analysis of some very limited examples. We still spent some time deploying mainly 2 core technologies that we detail hereafter.

7.1 DBpedia Spotlight

DBpedia Spotlight (Mendes et al. [5]) is an open-source free tool that allows to connect text to existing entries of Wikipedia and DBpedia. Basically, by giving the body of the article to the tool, we are able to extract the words and group of words that appear both in the article and as encyclopedic/ontological entries in those resources.

After running this tool we were able to extract 14 867 719 keywords from the 901 156 articles. This represents an average of 16.5 keywords per article. An excerpt of an article with the extracted keywords can be seen in Figure 6.

<p>Trois ingénieurs québécois sur quatre estiment que le gouvernement Charest doit adopter un moratoire complet sur l'exploration et l'exploitation du gaz de schiste en attendant le résultat d'études environnementales complètes, indique un sondage réalisé pour le compte du Réseau des ingénieurs du Québec. Selon le sondage, effectué par la firme Senergis, trois ingénieurs sur cinq sont pour l'heure défavorables à l'exploitation des gaz de schiste au Québec.</p> <p>[...]</p> <p>Trois répondants sur cinq sont aussi d'avis que le BAPE ne réussira pas à apporter des réponses satisfaisantes aux préoccupations de la population.</p>	<table> <tr> <th colspan="2">DBpedia Spotlight</th></tr> <tr> <td>Québec</td><td></td></tr> <tr> <td>gouvernement charest</td><td></td></tr> <tr> <td>gaz de schiste</td><td></td></tr> <tr> <td>BAPE</td><td></td></tr> <tr> <td>québécois</td><td></td></tr> </table>	DBpedia Spotlight		Québec		gouvernement charest		gaz de schiste		BAPE		québécois	
DBpedia Spotlight													
Québec													
gouvernement charest													
gaz de schiste													
BAPE													
québécois													

Figure 6: Excerpt of an article and extracted keywords from DBpedia Spotlight.

7.2 YAKE

Yet Another Keyword Extractor (YAKE) (Campos et al. , [1]) is also a free open-source tool to extract keywords. The main difference is that YAKE is completely language independent and instead of mapping words to encyclopedic entries, it is based on the word frequencies, their location in sentence, named entities indicators, etc.

This means that it extracts much more potential keywords but also much more noise as we can see in Figure 7.

7.3 Future work

Concerning the keyword extraction task, we already have results but we lack the evaluation tools to test them. We thought of multiple ways we could have done so, if we had more time.

One way would have been to run both these tools on keyword labeled data and compare the output result with the human labels. Our guess is that this method would have rendered a very low precision score, since the algorithms are not able to detect which keywords are representative of the whole article the same way a human annotator can.

Another evaluation strategy we could have used is the human evaluation. By detecting all the keywords from the keywords extraction tools and having them analyzed by a jury of annotators, we would be able to get an extremely precise way to know if the task is up to human standards. Nevertheless, this is extremely time and cost intensive

<p>Trois ingénieurs québécois sur quatre estiment que le gouvernement Charest doit adopter un moratoire complet sur l'exploration et l'exploitation du gaz de schiste en attendant le résultat d'études environnementales complètes, indique un sondage réalisé pour le compte du Réseau des ingénieurs du Québec. Selon le sondage, effectué par la firme Senergis, trois ingénieurs sur cinq sont pour l'heure défavorables à l'exploitation des gaz de schiste au Québec.</p> <p>[...]</p> <p>Trois répondants sur cinq sont aussi d'avis que le BAPE ne réussira pas à apporter des réponses satisfaisantes aux préoccupations de la population.</p>	YAKE
	ingénieurs
	gouvernement charest
	adopter un moratoire
	d'études environnementales complètes
	québec
	gaz de schiste
	sondage
	réseau

Figure 7: Excerpt of an article and extracted keywords from DBpedia Spotlight.

since we would need at the very least 3 human annotators reading each article and analyzing the keywords that were extracted, and on a sample set big enough to be representative.

Another choice would be to annotate our corpus with the keywords from each system and seeing if they improve the results on our theme/sub-theme classification task.

8 Conclusion

We had 2 tasks for this workshop: to classify articles by themes and sub-themes based on their content and to automatically extract the keywords. For the first task we proposed four different systems rendering very close results: one based on Google's BERT model, one based on Facebook's Fasttext model, one based on the logistic regression approach and one based on the SVM approach. The classification sub-tasks were evaluated by comparing the top 1 prediction of the model to the human annotated label. For the theme classification sub-task, the SVM returns the best result of all but just by a 3% difference. For the sub-theme classification problem, we experimented with only two out of four methods, of which, Fasttext rendered the best result.

The theme/sub-theme classifiers allow to assist the journalists and content editors and save them the time to search potential themes and sub-themes to classify each new article. The keyword extractor allows to add content metadata to the already existing article database. Both tasks could be further improved in many ways but this would require an even more profound analysis of the data annotation and an implementation of keyword extraction evaluation.

References

- [1] Campos, Ricardo and Mangaravite, Vítor and Pasquali, Arian and Jorge, Alípio Mário and Nunes, Célia and Jatowt, Adam. A text feature based automatic keyword extraction method for single documents. In *European Conference on Information Retrieval*, pages 684–691. Springer, 2018.

- [2] Cortes, Corinna and Vapnik, Vladimir. Support-vector networks. In *Machine learning*, volume 20, number 3, pages 273–297. Springer, 1995.
- [3] Devlin, Jacob and Chang, Ming-Wei and Lee, Kenton and Toutanova, Kristina. Bert: Pre-training of deep bidirectional transformers for language understanding. In *arXiv preprint*, arXiv:1810.04805, 2018.
- [4] Joulin, Armand and Grave, Edouard and Bojanowski, Piotr and Mikolov, Tomas. Bag of tricks for efficient text classification. In *arXiv preprint*, arXiv:1607.01759, 2016.
- [5] Mendes, Pablo N and Jakob, Max and García-Silva, Andrés and Bizer, Christian. DBpedia spotlight: shedding light on the web of documents. In *Proceedings of the 7th international conference on semantic systems*, pages 1–8. ACM, 2011.