# A deep learning classifier for sentence classification in biomedical and computer science abstracts

**Sérgio Gonçalves · Paulo Cortez · Sérgio Moro**

**Abstract** The automatic classification of abstract sentences into its main elements (background, objectives, methods, results, conclusions) is a key tool to support scientific database querying, to summarize relevant literature works and to assist in the writing of new abstracts. In this paper, we propose a novel deep learning approach based on a convolutional layer and a bi-directional gated recurrent unit to classify sentences of abstracts.
First, the proposed neural network was tested on a publicly available repository containing 20 thousand abstracts from the biomedical domain. Competitive results were achieved, with weight-averaged precision, recall and F1-score values around 91%, and an Area Under the ROC Curve (AUC) of 99%, which are higher when compared to a state-of-the-art neural network.
Then, a crowdsourcing approach using gamification was adopted to create a new comprehensive set of 4,111 classified sentences from the computer science domain, focused on social media abstracts. The results of applying the same deep learning modeling technique trained with 3,287 (80%) of the available sentences were below the ones obtained for the larger biomedical dataset, with weight-averaged precision, recall and F1-score values between 73% and 76%, and an AUC of 91%. Considering the dataset dimension as a likely important factor for such performance decrease, a data augmentation approach was further applied. This involved the use of text mining to translate sentences of the computer science abstract corpus while retaining the same meaning. Such approach resulted in slight improvements, of around 2 percentage points, for the weight-averaged recall and F1-score values.

S.Gonçalves and P. Cortez
ALGORITMI Centre, Dep. Information Systems, University of Minho, Portugal
E-mail: a72886@alunos.uminho.pt, pcortez@dsi.uminho.pt

S. Moro
Instituto Universitário de Lisboa (ISCTE-IUL), ISTAR-IUL, Lisboa, Portugal
E-mail: sergio.moro@iscte-iul.pt

## 1 Introduction

The number of scholarly works has increased during in the last decades [24]. For example, around 114 million of English scholarly documents were accessible on the Web in 2014 [18]. Such volume makes it difficult to quickly select relevant scientific documents. Scientific abstracts summarize the most important elements of a paper and thus those are valuable sources for filtering the most relevant papers during a literature review process [1].

The classification of scientific abstracts is a particular instance of the sequential classification task, considering there is a typical order in the classes (e.g., the 'Objective' label tends to appear after the 'Background') [50]. This classification transforms unstructured text into a more information manageable structure [11]. This is acknowledged by the Emerald publisher, which requires all submissions to include a structured abstract [7]. In effect, the automatic classification of abstract sentences presents several advantages. It is a valuable tool for general scientific database querying (e.g., using Web of Science, Scopus). Also, it can assist in manual [20] or text mining [26] systematic literature review processes, as well as other bibliometric analyses. Moreover, it can help in the writing of new paper abstracts [22].

In this study, we present a novel deep learning neural network architecture for the sequential classification of abstract sentences. The architecture uses a word embedding layer, a convolutional layer, a bi-directional Gated Recurrent Unit (GNU) and a final concatenation layer. The proposed deep learning model is compared with a recently proposed bi-directional Long Short-Term Memory (LSTM) based model [11], showing an interesting performance on classifying sentences under five classes: 'Background', 'Objectives', 'Methods', 'Results' and 'Conclusions' in abstracts corpus from two distinct domains. First, the approach is evaluated using a large publicly available 20K biomedical abstract corpus [10], which was also previously studied using a different approach [11], thus enabling a direct comparison. As a secondary contribution, we build a purely new abstract corpus from the computer science domain, with a particular focus on the social media topic. The latter dataset was created by adopting a crowdsourcing approach with gamification features to attract researchers to manually classify the sentences. By building on collective intelligence, more than 4k sentences were categorized. To further improve the classification results, we employed data augmentation based on text mining, thus doubling the training set size. The new corpus was also made publicly available, thus it can be used in future research comparison studies. We consider two distinct abstract sentence corpus, from different domains (biomedical and computer science), in order to provide a more robust validation. The good classification results that were obtained for both domains (Section 4) suggest that our ap-

proach is potentially valuable for sentence classification of abstracts from any scientific domain.

This paper is organized as follows: Section 2 presents the related work; Section 3 describes the two abstract corpus, the deep learning architecture and evaluation metrics; the obtained results are analyzed in Section 4; finally, conclusions are drawn in Section 5.

## 2 Related Work

### 2.1 Deep Learning for Abstract Sentence Classification

As pointed out by Dernoncourt et al. [11], most sequential sentence classification methods are based on 'shallow' methods (e.g., naive Bayes, Support Vector Machines (SVM)) that require a manual feature engineering based on lexical (e.g., bag of words, n-grams), semantic (e.g, synonyms), structural (e.g., part-of-speech tags) or sequential (e.g., sentence position) information. The advantage of using deep learning is that the neural networks do not require such manual design of features. Also, deep learning often achieves competitive results in text classification [16].

Regarding abstract sentence classification, this topic has been scarcely researched when compared to other text classification tasks (e.g., sentiment analysis). The main reason for this reduced attention is the restricted availability of public datasets. In 2010 [3], the manual engineering approach was used to set nine features (e.g., bi-grams) and train five classifiers (e.g., SVM) that were combined to classify four main elements of medical abstracts. In 2013 [22], a private corpus with 4550 abstracts from different scientific fields was collected from ScienceDirect. The abstract sentences were manually labeled into four categories: 'Background', 'Goal', 'Methods' and 'Results'. The authors also used the conventional manual feature design approach (e.g., n-grams) and a transductive SVM. More recently, in 2017 [10], a large abstract corpus was made publicly available. Using this dataset, a deep learning model, based on one bi-directional LSTM, was proposed for a five class sentence prediction, outperforming four other approaches (e.g., n-gram logistic regression, multilayer perceptron) [11]. In this paper, we propose a different deep learning architecture, mainly composed by a convolutional layer and a bi-directional GRU layer to classify the sentences from abstracts, which uses word embeddings instead of character embeddings. By taking into consideration the position of the sentences, as well as encoding contextual information on the vector of each sentence,we expect that the proposed architecture can potentially achieve better results when compared to the study by Dernoncourt et al. [11].

### 2.2 Crowdsourcing and Collective Intelligence to Classify Data

The difficulties of gathering already classified datasets highlighted in the previous subsection raises the challenge of testing our approach in several domains.

Yet, it is essential for a broader validation of the accuracy of our proposal. Furthermore, by building a new dataset and making it publicly available, it can constitute a baseline for future comparisons with distinct approaches of ours, helping to pave avenues for future research. Thus, this study sets the additional challenge of gathering a newly classified dataset from another domain to complement the biomedical dataset made publicly available by Dernoncourt and Lee [10].

Collective intelligence consists in benefiting from many individuals by combining their intelligence in cognitive tasks [9, 43]. Currently, the Internet has leveraged collective intelligence by bringing communities together online in discussion forums and groups in social networks. Thus, the technology currently available at the reach of a click instantaneously connects people who may be spread across the globe. As a result, the concept of crowdsourcing has emerged to take advantage from collective intelligence by using the Internet to gather people together toward a common goal of solving cognitive tasks [4], from data classification to solving complex mathematics problems. Specifically, Welinder and Perona [44] have shown that data annotation can be effectively done using a web platform available online offering a crowdsourcing service.

One of the main challenges of crowdsourcing is to keep participants motivated during the whole process. Thus, a crowdsourcing platform must include an incentive system to prompt for participation through positive encouragement [29]. Incentives and motivation are interconnected, since the former helps in building the latter [17]. Therefore, adopting strategies that keep individuals motivated is at the core of crowdsourcing [23]. The motivation may be extrinsic, through external incentives such as money or discount coupons, or intrinsic, when the individual is pleased just by participating on crowdsourcing [17]. While both extrinsic and intrinsic motives are shown to influence participation, some authors such as Zheng et al. [49] argue that the latter motives can be more effective than the former, although such difference may be also due to cultural factors. A possible solution freely available of intrinsic incentive is gamification. It offers users appealing awards granted on desired accomplishments to incentivize participation [28].

## 3 Materials and Methods

### 3.1 Computer Science Abstract Corpus Retrieval and Classification

The first task to build the comprehensive set of classified abstract sentences in computer science applied to social media is to retrieve the raw set of abstracts. Although there are several platforms which freely publish scientific abstracts, most of them do not show a well-defined policy about what can be publicly shown to third-parties and, specifically, in the form of a dataset. As such, we choose the arXiv, which allows using articles' metadata, including abstracts, and provides an API for both commercial and non-commercial usage[1]. To re-

---

[1]  https://arxiv.org/help/oa/index

trieve data, we adopted the *aRxiv R* package, which uses the API to efficiently search the archive. Within the "computer science" category, we searched only for articles containing the terms "social media", "social network" or "social networks" in the title.

The different human languages pose serious challenges in text mining and machine learning tasks and, as such, many researchers choose to focus in one language [e.g., 5]. Since most of scientific literature is currently published in the English language [12] and, additionally, the large publishers and top journals require submissions to be in this language, we filtered all the collected articles using the *textcat* package [14] to match English abstracts. The result is that from a total of 658, only 4 were written in other than the English language.

We implemented our crowdsourcing solution using the Amazon Elastic Compute Cloud, which is an Infrastructure–as–a–Service cloud platform. This is the platform with the largest market share, with around 40% of the Infrastructure–as–a–Service market [32], and it is freely available for academic purposes. The Ubuntu 16.04 operative system and the MariaDB database system were installed in the cloud, as well as the Nginx web server, which presents better performance and less memory requirements when compared to Apache, according to Suciu et al. [41]. We adopted the PHP Laravel (PHP 7.2) framework for developing our solution.

We made available our web responsive platform (adaptable depending on the type of device) in the following URL: http://classifyabstracts.info/. It was shared in academic social networks such as "www.researchgate.net", as well as through our own research contacts. The user needs to login to access an initial small tutorial that introduces the aims and scope of this research. This introductory tutorial consists in two steps. The first consists in reading a few examples of sentences previously classified, such that the user takes contact with the possible categories for the classification procedure (Figure 1).

On the second and last step, the user needs to correctly classify each sentence of an abstract already classified. To simplify the annotation process, the categories for classification are selected through drop–down lists. If the user has at least one wrong classification, a warning appears highlighting the mistake and the user needs to repeat this step again. After successfully completing the tutorial, the user may start classifying unclassified abstracts in the platform (Figure 2).

After the user classifies an abstract, a message is shown displaying the number of already classified abstracts. As an incentive, a simple gamification component consisting in points (one point per classified abstract) and three levels was introduced. Thus, a user with 10 points achieves the bronze level, while a user reaching 50 points achieves the silver level and, finally, a user with 100 or more points achieves the gold level.

In total, 76 users classified abstracts, with very few of them classifying more than 100 abstracts, and still few reaching the silver level. Thus, most users devoted little effort and classified few abstracts. This is a known inequality

**Fig. 1** First step of the introductory tutorial
.

phenomenon of crowdsourcing participants, in which only a small fraction of users holds a large contribution [39].

### 3.2 Scientific Abstract Corpora

Data understanding is a critical step in the knowledge discovery process [25]. Thus, in this subsection we characterize each of the two datasets used to test our approach. One of them is the one gathered through crowdsourcing (Section 3.1), while the other is the PubMed 20k dataset made publicly available by Dernoncourt and Lee [10]. The latter, from here forth referred to as "PubMed 20k", sets the baseline for comparison purposes [46]. The former, labelled as "CS Abstracts", enables to test our approach in a new domain, and additionally provides to future researchers another corpus for sentence classification. The corpus of the PubMed 20k includes open access papers from the PubMed biomedical database and related with Randomized Controlled Trials (RCT).

**Fig. 2** Second step of the introductory tutorial.

The sentences were classified by the authors of the articles into the five stan-dardized labels. The full corpus has a total of 200K abstracts. A smaller subset, with 20K most recent abstracts, was also made available for a faster experi-mentation of sequential sentence classification methods. Considering the 20K subset was used in the work of Dernoncourt et al. [11], we also adopted the same dataset, to facilitate the experimental comparison.

The size of both datasets is substantially different, with the PubMed 20k containing a larger number of abstracts and, subsequently, more classified sentences. Table 1 shows the total number of abstracts and sentences used for training, validating the model, and testing it. The abstract distribution through the three partitions in the "CS Abstracts", training, validation and testing, was done to evenly maintain distribution of each class. According to Ng [31], a uniform distribution between the three partitions in the datasets is used to prevent unexpected performance issues in productive environments.

**Table 1** Number of abstracts and sentences for each dataset by usage (i.e., training, validating and testing).

| Dataset | | Training | Validation | Testing |
|---|---|---|---|---|
| PubMed 20k | Nr. abstracts | 15,000 | 2,500 | 2,500 |
| | Nr. sentences | 180,000 | 30,000 | 30,000 |
| CS Abstracts | Nr. abstracts | 500 | 77 | 77 |
| | Nr. sentences | 3,287 | 824 | 619 |

A model is trained using the training dataset, it is tuned according to the obtained results in the validation dataset and its generalization capability is measured on the test dataset. As for the PubMed 20k dataset, it is made publicly available with the data already partitioned as shown in Table 1. Figure 3 highlights the different distributions of the two datasets. Such result can be justified by the different data collection approaches and subsequent writing styles of both sciences, biomedical and computer science. Myers et al. [30] corroborate our claim by highlighting that the writing style is interconnected to the research method. This result supports the importance of our crowdsourcing approach in gathering a new dataset from another domain to validate our classifier, which is emphasized by the lack of publicly available classified datasets. Figure 3 also confirms that the different train, validation and test partitions have similar class distributions, which is a desired trait for classification tasks [31].

### 3.3 Neural Networks Models

In the last years, there has been remarkable developments in deep learning [16][47]. Architectures such as Convolutional Neural Network (CNN), LSTM and GRU have obtained competitive results in several competitions (e.g., computer vision, signal and natural language processing).

The CNN is a network mainly composed by convolutional layers. The purpose of the convolutional layers is to extract features that preserve relevant information from the inputs [21]. To obtain the features, a convolutional layer receives a matrix as input, to which a matrix with a set of weights, known as a filter, is applied using a sliding window approach and, at each of the sliding window steps, a convolution is calculated, resulting in a feature. The size of the filter is a relevant hyperparameter.

Although CNNs have been widely used in computer vision, they can also be used in sentence classification [19]. The use of convolutional layers enables the extraction of features from a window of words, which is useful because word embeddings alone are not able to detect specific nuances, such as double negation, which is important for sentiment classification. The width of the filter, represented by $h$, determines the length of the n-grams. The number of filters is also a hyperparameter, making it possible to use multiple filters with

**Fig. 3** Relative class frequencies for the CS Abstracts (top bar charts) and PubMed 20k (bottom bar chars).

varying lengths [19]. The filters are initialized with random weights and, during network training, weights are learned for the specific task of the network, through backpropagation. Since each filter produces its own feature map, there is a need to reduce the dimensionality caused by using multiple filters. A sentence can be encoded as a single vector by applying a max pooling layer after the convolutional layer, which takes the maximum value for each position, from all the feature maps, keeping only the most important features.

Recurrent Neural Networks (RNN) are relevant for sequential data, such as the words that appear in a sentence. Consider the words $(x_1, ..., x_t)$ from a

given sentence (sequence of words). The hidden state $s_t$ of the word $x_t$ depends on the hidden state $s_{t-1}$, which in turn is the hidden state of the word $x_{t-1}$ and, for this reason, the order in which words appear over the sequence also influences the various hidden states of the RNN.

The LSTM network is a particular RNN that uses an internal memory to keep information between distant time steps to model long-term dependencies of the sequence. It uses three gating mechanisms, input gate, forget gate and output gate, which control (at each hidden state), what new information should be updated into the memory, and what information should be erased from the memory. The GRU [6] was recently introduced and it can be used as an alternative to the LSTM model. The GRU uses a reset and update gates, which are able to control how much information should be kept from previous time steps. Both GRU and LSTM are solutions that help mitigate the vanishing gradient problem of conventional RNNs.

A deep learning model was used by Dernoncourt et al. [11] for abstract sentence classification. The model uses character embeddings that are then concatenated with word embeddings and used as input for a bi-directional LSTM layer, which outputs a sentence vector based on those hybrid embeddings. The sentence vector is used to predict the probabilities of the labels for that sentence. The authors also use a sequence optimization layer, which has the objective of optimizing the classification of a sequence of sentences, exploiting existing dependencies between labels.

## 3.4 Proposed Architecture

The proposed word embedding, convolutional and bi-directional GRU (Word-BiGRU) architecture is shown in Figure 4. We assume that each abstract has $i$ sentences $(S_1, ..., S_i)$ and each individual sentence has $n$ words $(x_1^1, ..., x_n^i)$, where $x_n^i$ is the $n^{th}$ word from the $i^{th}$ sentence. The various words from the sentences are mapped to their respective word embeddings, and those embedding are used to create a sentence matrix $E \in R^{m \times d}$, where $d$ equals to the dimensionality of the embeddings. We use word embeddings pre-trained on English Wikipedia, provided by Glove (with $d = 200$) [33].

Then, a convolutional layer is used with a sliding window approach that extracts the most important features from the sentences. Let $E \in R^{m \times d}$ denote the sentence matrix, $w \in R^{h \times d}$ a filter, and $E[i:j]$ the sub-matrix from row $i$ to $j$. The single feature $o_i$ is obtained using:

$$o_i = w * E[i : i + h - 1] . \tag{1}$$

In this study, we use a filter with a size of $h = 5$. To add nonlinearity to the output, an activation function applied to every single feature. For the feature $o_i$, it is obtained by:

$$c_i = f(o_i + b); \tag{2}$$

**Fig. 4** Schematic of the proposed Word-BiGRU deep learning architecture.

where $f$ is the activation function and $b$ is the bias. We use ReLU as the activation function in our model because it tends to present a faster convergence [15].

Next, we take the various features maps obtained from the convolutional layer, and feed them into a max pooling layer to encode the most important features extracted by the convolutional layer into a single vector representation that can be used by the next layers. Let $g_1, ..., g_i$ denote several vectors, each one encoding a particular sentence of the abstract. The vectors are then fed to bi-directional GRU layer, where the hidden states for each time step are calculated. We will use $\odot$ to denote the Hadamard Product, while using $W$ and $U$ to denote weight matrices of the GRU layer. Let $h_{i-1}$ be the hidden state of the previous sentence from the same abstract, the candidate hidden state $\tilde{h}_i$ for the current sentence is given by:

$$\tilde{h}_i = \tanh(W_h g_i + U_h(r_i \odot h_{i-1}) + b_h) \, . \tag{3}$$

The reset gate $r_i \in [0, 1]$ has the purpose of controlling how much information of the past hidden state, $h_{t-1}$ will be kept. Let $\sigma$ be the activation function, the reset gate $r_i$ is calculated by:

$$r_i = \sigma(W_r g_i + U_r h_{i-1} + b_r) \, . \tag{4}$$

To control how much new information will be stored in the hidden state, an update gate $z_i \in [0,1]$ is used, given by:

$$z_i = \sigma(W_z g_i + U_z h_{i-1} + b_z) \ . \tag{5}$$

The hidden state $h_i$, which is the hidden state of the sentence $i$, is obtained by:

$$h_i = z_i \odot \tilde{h}_i + (1 - z_i) \odot h_{i-1} \ . \tag{6}$$

Since we use a bi-directional GRU layer, there is a forward pass and a backward pass. The hidden states resulting from the forward pass are:

$$(\overrightarrow{h_1}, ..., \overrightarrow{h_i}) \ . \tag{7}$$

where $h_i$ is the hidden state of the $i^{th}$ sentence of the abstract. Similarly, the hidden states resulting from the backward pass are:

$$(\overleftarrow{h_1}, ..., \overleftarrow{h_i}) \ . \tag{8}$$

By using a bi-directional GRU, we want to capture contextual information about each sentence of the abstract, by taking into consideration the sentences that appear before and after it. For the $i^{th}$ sentence of the abstract, the individual vector $k_i$, which encodes the sentence with contextual information captured using the bi-directional GRU layer, is obtained by concatenating the forward and backward hidden states:

$$k_i = [\overrightarrow{h_i} \oplus \overleftarrow{h_i}]; \tag{9}$$

where $\oplus$ is the concatenation operator. Each encoded sentence $k_i$ is then concatenated with an integer value indicating the position of that sentence in the abstract, resulting in $z_i$:

$$z_i = [k_i \oplus i] \ . \tag{10}$$

Finally, a softmax layer is used, such that the outputs can be interpreted as class probabilities.

## 3.5 Evaluation

In multiclass tasks, a classifier often outputs a class probability and the highest probability class is assigned as the predicted class label. Using these predicted labels, classification accuracy is often measured by building a confusion matrix, which maps predicted versus desired labels. From this matrix, several metrics can be computed, such as [45]: Precision, Recall, F1-score. For a class $c$, these metrics are obtained using:

$$\begin{aligned} Precision_c &= \frac{TP_c}{TP_c + FP_c} \\ Recall_c &= \frac{TP_c}{TP_c + FN_c} \\ F1 - score_c &= 2 \times \frac{Precision_c \times Recall_c}{Precision_c + Recall_c} \ . \end{aligned} \tag{11}$$

where $TP_c$, $FP_c$, $FN_c$ denote the number of true positives, false positives and false negatives for class $c$.

Another possibility to analyze multiclass probabilities is to consider one class probability $p_c$ and a decision threshold $K \in [0, 1]$. The class is considered positive if $p_k > K$. The receiver operating characteristic (ROC) curve shows the performance of the classifier across all $K$ values for class $c$, plotting one minus the specificity ($x$-axis) versus the sensitivity ($y$-axis) [27]. The overall discriminatory performance is given by the area under the curve ($AUC = \int_0^1 ROCdK$). It is common to interpret the quality of the AUC values as: 0.5 – equal to a random classifier; 0.6 – reasonable; 0.7 – good; 0.8 – very good; 0.9 – excellent; and 1 – perfect.

To combine all five class confusion matrices results into a single measure, we adopt two aggregation methods: macro-averaging and weight-averaging. The macro-averaging computes first the metric (e.g., Precision using Equation 11) for each class and then averages the overall result. The weight-averaging is computed in a similar way except that each class metric is weighted proportionally to its prevalence in the data. Dernoncourt et al. [11] used only the weight-averaging method. As for the ROC analysis, the five AUC class values are aggregated by using a macro-average aggregation method [13].

For comparison purposes, we adopt the same train, validation and test sets used by Dernoncourt et al. [11] for the PubMed 20k corpus. The respective set partition numbers are shown in Table 1. This table also includes the new CS abstract corpus training, validation and test division adopted numbers that can be used in future comparison works.

## 4 Results

### 4.1 Hyperparameter Tuning

A predictive model needs to be tuned through an array of possible hyperparameters which may have a positive or negative impact in the model's predictive performance [40]. However, finding the best combination of hyperparameters which optimizes a model's performance is a complex task, given the possible combinations to be tested and the required time to train a model. Also, some hyperparameters are encompassed within a numeric interval, which makes it impossible to test all values. In this section, we report the effort made in tuning the Word-BiGRU model using the PubMed 20k dataset, which is a larger dataset for which there are previous baseline results provided by Dernoncourt et al. [11].

The Word-BiGRU model shares some components with the CNN, such as the convolutional layer structure. Nevertheless, the Word-BiGRU model training takes significantly more time when compared to the CNN. Thus, the hyperparameters optimization for the convolutional layer was done using the CNN model to speed up the process. Therefore, the goal is to find the best

values for the hyperparameters of the CNN model and then use those values in the convolutional layer of the Word-BiGRU model.

The filter size, number of filters, and dropout are the main hyperparameters of the CNN model. The filter size is accountable for finding how many words within a sentence are encompassed by the filter in the sliding window process for the convolution operation. The number of filters represents how many filters are used in the convolutional layer, with each filter being capable of extracting different features. The dropout is used to reduce overfitting to training data [38]. Thus, the following configuration was considered as a baseline for the CNN model: filter size – 32; number of filters – 3; dropout – 0.1. Also, to limit the space of possible solutions, the hyperparameter values were searched within the following intervals: filter size $\in \{3, 5, 8\}$; number of filters $\in \{32, 64, 128, 256\}$; dropout $\in \{0.1, 0.35, 0.5\}$.

Using such baseline, we explored the combinations of possible values. First, we changed the values of filter size, followed by the number of filters and, finally, by the dropout hyperparameter. Thus, this means that by exploring the number of filters, the best filter size value is already applied, since it was the first to be explored. According to Ng [31], hyperparameter selection should be done using one selected metric measured on the validation set. We follow such approach, in which we selected the weight-averaging method to aggregate the individual class metrics. The selected hyperparameter values are shown in Table 2. Thus, the best results for CNN within the defined intervals were obtained with 128 filters, each with size equals to 5, using a dropout of 0.35 after the convolutional layer.

**Table 2** CNN model hyperparameters (in %, validation set results, best values in **bold**).

| Hyperparameter | Valor | Precision | Sensitivity | F1-score |
|---|---|---|---|---|
| Number of filters | 32 | 82.0 | 81.9 | 81.9 |
| | 64 | 82.6 | 82.3 | 82.4 |
| | 128 | **83.3** | **83.4** | **83.3** |
| | 256 | 83.2 | 83.2 | 83.2 |
| Filter size | 3 | 83.3 | 83.4 | 83.3 |
| | 5 | **83.5** | **83.5** | **83.5** |
| | 8 | 82.9 | 82.2 | 82.5 |
| Dropout | 0.1 | 83.5 | 83.5 | 83.5 |
| | 0.35 | **83.6** | **83.6** | **83.6** |
| | 0.5 | 81.5 | 80.6 | 81.0 |

For the specific GRU bidirectional layer of Word-BiGRU, we assumed the previously selected CNN hyperparameters and ranged the number of GRU units within $\{25, 50, 75, 100\}$, as shown in Table 3. Thus, the selected number of GRU units was 50. The final adopted configuration of the Word-BiGRU model, used for both PubMed 20k and CS Abstracts corpora, includes: number of filters – 128; filter size – 5; dropout – 0.35; and number of GRU units – 50.

**Table 3** Word-BiGRU model hyperparameters (in %, validation set results, best results in **bold**).

| Hyperparameter | Value | Precision | Sensitivity | F1-score |
|---|---|---|---|---|
| Number of units | 25 | 89.4 | 89.3 | 89.3 |
| | 50 | **91.5** | **91.5** | **91.5** |
| | 75 | 89.8 | 89.7 | 89.7 |
| | 100 | 90.0 | 90.0 | 90.0 |

## 4.2 Biomedical Abstracts

The performance in the PubMed 20k dataset can be compared with the results achieved by Dernoncourt et al. [11] (Char-BiLSTM). The latter uses word characters as inputs of a bidirectional LSTM layer, which provides vectorial representations of those words. Each vector is then used as input of another bidirectional LSTM layer which, in turn, provides a vectorial representation of a sentence. Additionally, Dernoncourt et al. [11] defined a transition matrix which contains the probabilities of transition between classes. The values of this matrix are computed through the model which, given a sentence being classified to a class (e.g., "Objective"), can then compute the probability of the next sentence belonging to another class (e.g., "Methods"). The results of our approach (Word-BiGRU) are also compared to the standard CNN model (described in Section 3.3) in Table 4. The best results for each metric are shown in bold. Our approach (Word-BiGRU) outperforms the alternatives for all computed metrics.

**Table 4** Averaged test results for PubMed 20k (in %, best values in **bold**).

| Metric | Averaged | Char-BiLSTM [11] | CNN | Word-BiGRU |
|---|---|---|---|---|
| Precision | Macro-Averaged | 86.4 | 80.7 | **86.7** |
| | Weight-Averaged | 90.1 | 83.6 | **90.9** |
| Recall | Macro-Averaged | 83.7 | 77.6 | **86.7** |
| | Weight-Averaged | 89.9 | 83.5 | **90.8** |
| F1-score | Macro-Averaged | 85.0 | 78.5 | **86.7** |
| | Weight-Averaged | 90.0 | 83.5 | **90.8** |

The obtained results provide evidence that including contextual information improves performance for the PubMed 20k dataset. This is the main difference between the CNN and the Word-BiGRU, with the latter including contextual information through the GRU bidirectional layer. Such difference had a significant impact in all metrics.

The Word-BiGRU systematically outperforms the Char-BiLSTM model in all metrics, with improvements ranging from 0.3 (Precision Macro-Averaged) to 1.7 (F1-score Macro-Averaged) percentage points. Figure 5 shows the ROC curves, and respective AUC values (area), for each class. The average value

computed through macro–average is 0.99. This is very close to 1, suggesting
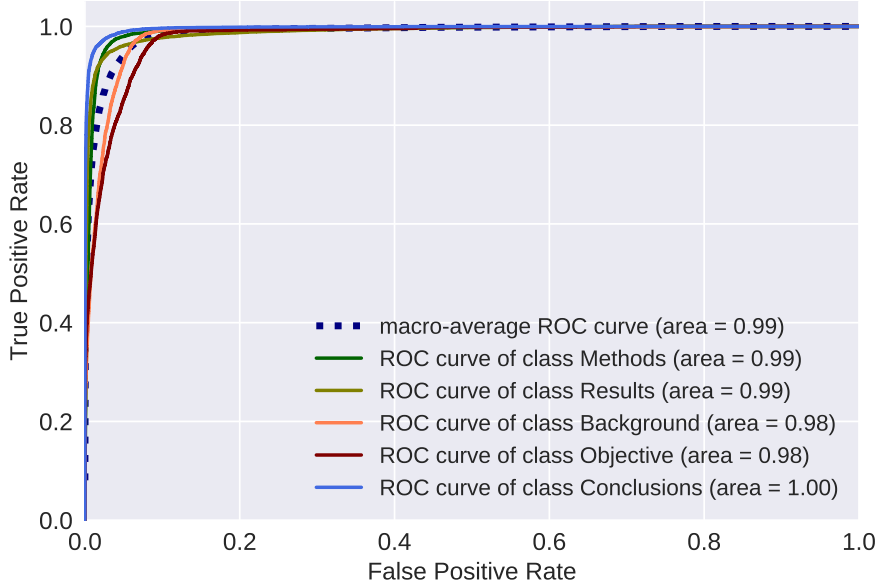the model has an excellent discriminatory capability.



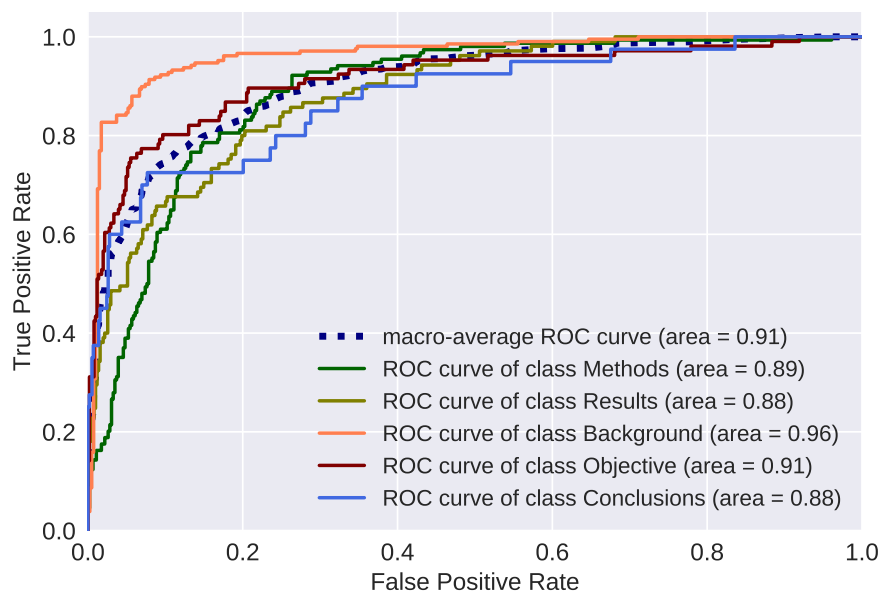**Fig. 5** ROC curves and AUCs for the Word-BiGRU model in the PubMed 20k dataset.

4.3 Computer Science Abstracts

Table 5 shows the classification performance test metrics for our Word-BiGRU
approach in comparison to the CNN model for the CS Abstracts corpus. The
World-BiGRU resuls are clearly better when compared with the CNN model,
meaning that including contextual information through the bidirectional GRU
layer had a significant impact on the model, confirming the results also ob-
tained for the PubMed 20k dataset. When comparing the Word-BiGru re-
sults for both datasets, lower metric values were achieved for CS Abstracts
(around 75%) when compared with PubMed 20k (around 91%). The classi-
fication performance differences can be due to several factors. For instance,
PubMed 20k sentences were classified by the own authors, while the CS Ab-
stracts were labeled using anonymous volunteers, via crowdsourcing. Further-
more, the PubMed 20k is around five times larger in size when compared to
the CS Abstracts dataset and often deep learning models improve results with
big data. This last issue is further explored in Section 4.4.

**Table 5** Averaged test results for CS Abstracts (in %, best values in **bold**).

| Metric | Method | CNN | Word-BiGRU |
|---|---|---|---|
| Precision | Macro–Average | 61.3 | **75.4** |
| | Weight–Averaged | 66.3 | **76.3** |
| Recall | Macro–Averaged | 57.0 | **64.7** |
| | Weight–Averaged | 67.7 | **73.0** |
| F1-score | Macro–Averaged | 59.1 | **69.6** |
| | Weight–Averaged | 67.0 | **74.6** |

The ROC curves shown on Figure 6 highlight a different picture of Word-BiGRU results when compared to the same model applied to the PubMeds 20k dataset (Figure 5). The "Conclusions" class is the one with the lowest AUC (although still very good value of 88%), whereas in the PubMed 20k was the one with the highest AUC. Such result may derive from the different class distribution of both datasets (as shown in Figure 3).



**Fig. 6** ROC curves and AUCs for the Word-BiGRU model in the CS Abstracts dataset.

4.4 Data Augmentation of CS Abstracts

The CS Abstracts corpus has led to worse classification results when compared to the PubMed 20k data. Deep learning usually relies on large datasets for better training a model [34]. Thus, we argue that the CS Abstracts reduced size of 3,287 sentences for training the model poses a limitation that we intend to overcome. A possibility that has recently been explored is data augmentation, which aims to artificially increase the size of the training dataset [2]. In domains such as computational visualization, simple transformations such as translation or scale changes can be directly applied to the training images to create different images that hold the same information (i.e., classified under the same category) [35].

However, in text mining, simple transformations are not a viable solution to data augmentation since textual data has a specific nature in which changing the form of a word may render a sentence to have a completely different meaning. A possible solution to overcome this issue is to replace words within sentences by synonyms. For example, a list of synonyms implemented in some NLP tools may be used, as demonstrated by Zhang et al. [48]. An alternative is to adopt translation to compute sentences with the same meaning, as shown by Pavel Ostyakov on GitHub[2]. The procedure starts by translating a sentence to other language and then back to the original language. This approach has the advantage of changing several word combinations at the same time and word order within a sentence, while keeping the same meaning, helping to build a more diverse training corpus. As such, we adopted this latter approach using the Textblob Python package. Also, we chose Spanish as the in-between translation language since the adopted NLP package was validated with success using the same language [42]. Table 6 shows three CS Abstracts sentences in both the original and the transformed formats. It shows that the sentences are slightly changed while retaining the same meaning.

**Table 6** Comparison of original versus transformed sentences as a result of data augmentation

| Original sentence | Transformed sentence |
| --- | --- |
| *"Early detection of such compromised accounts is very important in order to control the damage."* [37] | *The early detection of such compromised accounts is very important to control the damage.* |
| *"In this work we propose a novel general framework for discovering compromised accounts by utilizing statistical text analysis."* [37] | *In this paper we propose a new general framework for discovering compromised accounts through the use of statistical text analysis.* |
| *"These are the anomalies caused by a user because of his/her variable behaviour towards different sources."* [37] | *These are the anomalies caused by a user due to their variable behavior towards different sources.* |

---

[2] https://github.com/PavelOstyakov/toxic/blob/master/tools/extend_dataset.py

This data augmentation process was only applied to the training set, which doubled in size from 3,287 to 6,574 sentences (i.e., for each sentence an equivalent one with the same meaning was added to the dataset). Table 7 shows the results of training the Word-BiGRU model with both the original and the data augmented datasets.

**Table 7** Comparative results between original and data augmented datasets for the Word-BiGRU model.

| Metric | Method | Original | Data augmented |
|--------|--------|----------|----------------|
| Precision | Macro–Averaged | **75.4** | 69.7 |
|  | Weight–Averaged | **76.3** | 75.4 |
| Recall | Macro–Averaged | 64.7 | **70.0** |
|  | Weight–Averaged | 73.0 | **74.9** |
| F1-score | Macro–Averaged | 69.6 | **69.8** |
|  | Weight–Averaged | 74.6 | **75.1** |

Using an augmented training set, when compared with the original dataset, leads to a decrease in Precision but an improvement in both Recall and F1-score measures. As for the ROC curves (Figure 7), the obtained results are quite similar, with the augmented trained model presenting slight AUC improvements for the Methods and Results classes, with percentage point differences of: 1 – "Methods" and 2 – "Results".

### 4.5 Sensitivity Analysis

All the models evaluated are based on deep learning, which are considered black-box models. One way to open these black-box models is to use a sensitivity analysis [8], which monitors the output responses when varying the input features through their range of possible values. Specifically for textual contents, Ribeiro et al. [36] have developed the Lime package in Python, which changes model's input by adding and removing words to assess the impact on the model's output. We have adopted this package to obtain the words in sentences that influence the most each category. The aim is to understand the differences in the decision making process of our approach (Word-BiGRU) in comparison to CNN.

For demonstration purposes, in this section we compare the sensitivity analysis of CNN and Word-BiGRU models when trained with the augmented CS Abstracts dataset. Figure 8 shows the obtained output, when using the Lime package, for one selected abstract. Words shaded in red are the ones that contributed the most to classifying the sentence as "Background", with particular emphasis for "media" and "development". The words shaded in green are those that provided additional hints that the sentence may belong
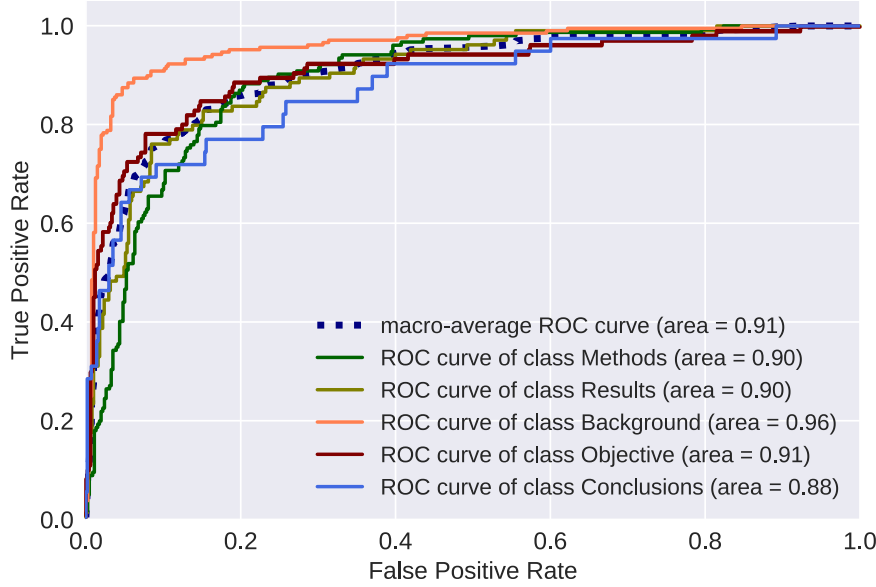
**Fig. 7** ROC curves and AUCs for the Word-BiGRU model in the original and data augmented CS Abstracts datasets.

to other class than "Background". The words shaded that belong to other sentences emerge because the Lime package affects the whole abstract.
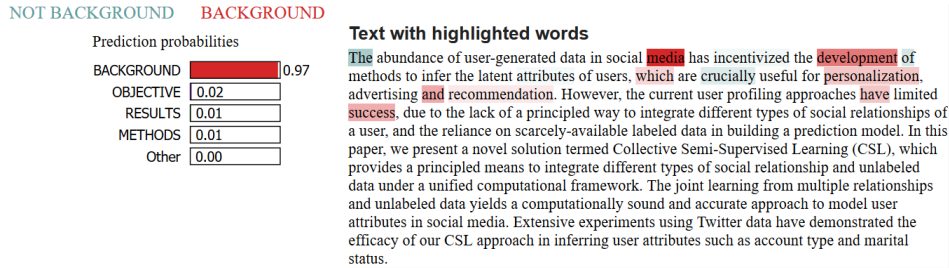


**Fig. 8** Most relevant words in classifying the first sentence for CNN model.

Figure 9 shows the Lime results for the Word-BiGRU model on the same selected abstract. The model has a stronger confidence that the first sentence is related with "Background". And some of the highlighted words for reaching to such conclusion are different (e.g., "latent attributes" instead of "development"). Another interesting feature of Word-BiGRU is that it uses several words from other sentences, as shown by looking at the shaded words in red

in the sentences following the first. This is an evidence of the contextual information used for feeding the model.
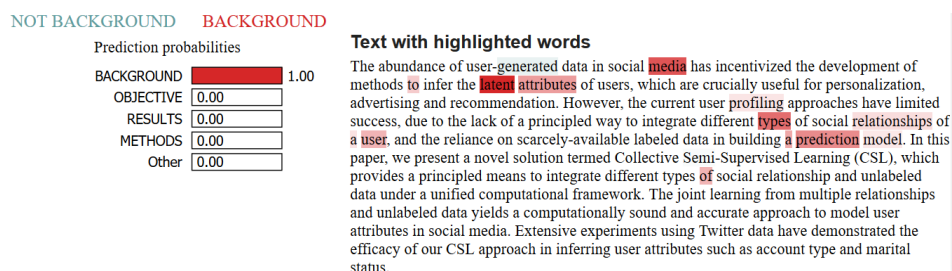


**Fig. 9** Most relevant words in classifying the first sentence for Word-BiGRU model.

The results of applying the Lime package to the CNN model for the second sentence of the abstract are shown on Figure 10. Again, the sentence is accurately classified as "Background". The words "however" and "current" were the ones influencing the most such classification.
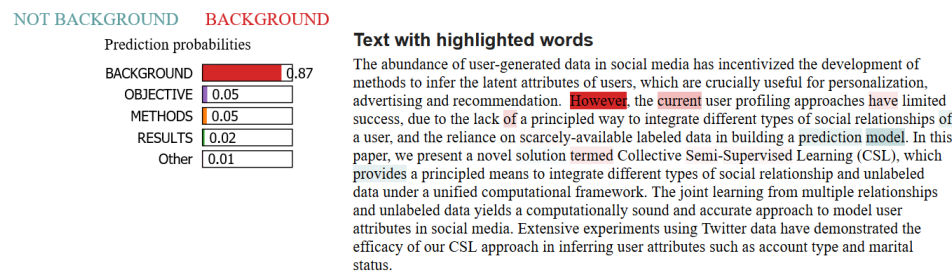


**Fig. 10** Most relevant words in classifying the second sentence for CNN model.

For the Word-BiGRU model, the probability of the second sentence being classified as "Background" is higher, again proving the value of the contextual information included in the model (Figure 11). This can be shown in the highlighted words in sentences other than the second.

## 5 Conclusions

This paper presents a novel deep learning architecture for the classification of scientific abstract sentences (background, objectives, methods, results, conclusions), which is valuable to assist in scientific database querying, performing literature reviews and to support the writing of new abstracts. The proposed Word-BiGRU architecture assumes word embeddings, a convolutional layer and a bi-directional Gated Recurrent Unit (GRU). Using a large sentence corpus, related with 20k abstracts from the biomedical domain, we have obtained
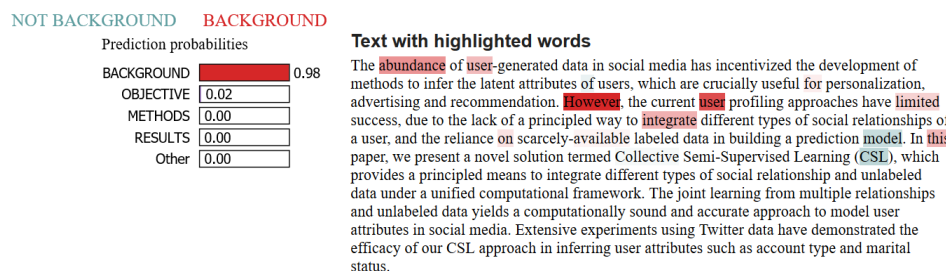
**Fig. 11** Most relevant words in classifying the second sentence for Word-BiGRU model.

high quality classification performances, with weight-averaged Precision, Recall and F1-score values around 91%. These results compare favourably against a state-of-the-art bi-directional Long Short-Term Memory (LSTM) model.

We also collected a new dataset of Computer Science scientific abstracts using a crowdsourcing approach to assess the performance of our architecture in a new domain. The results in this 4k classified sentences dataset were below the ones achieved for the 20k biomedical dataset when using the same Word-BiGRU architecture. However, a very good classification level was achieved by Word-BiGRU, which compared favorably with a Convolutional Neural Network (CNN) model, obtaining weight-averaged Precision, Recall and F1-score values around 75%. Furthermore, we confirmed the decision process of Word-BiGRU and CNN through a sensitivity analysis procedure, and we found evidences of the former using the contextual information included through the GRU. To address the reduced size of the collected Computer Science abstracts dataset, we adopted a data augmentation approach based on sentence translation to Spanish and back to English, which resulted in a training corpus with twice the original size. The classification metrics on test set data have shown a slight improvement of the augmented trained Word-BiGru model in terms of sensitivity and F1-scores.

In future research, we intend to check if we can improve the Computer Science corpus results by further increasing the Computer Science corpus size, either by collecting more human labels or by combining our data augmentation translation technique with other synthetic data creation methods (e.g., usage of synonyms). We also wish to enlarge the experimentation of the proposed Word-BiGRU deep learning architecture to other sequential tasks.

**References**

1. Atanassova I, Bertin M, Larivière V (2016) On the composition of scientific abstracts. Journal of Documentation 72(4):636–647

2. Bawa VS, Kumar V (2018) Emotional sentiment analysis for a group of people based on transfer learning with a multi-modal system. Neural Computing and Applications pp 1–12

3. Boudin F, Nie JY, Bartlett JC, Grad R, Pluye P, Dawes M (2010) Combining classifiers for robust pico element detection. BMC medical informatics and decision making 10(1):29

4. Brabham DC (2008) Crowdsourcing as a model for problem solving: An introduction and cases. Convergence 14(1):75–90

5. Canito J, Ramos P, Moro S, Rita P (2018) Unfolding the relations between companies and technologies under the big data umbrella. Computers in Industry 99:1–8

6. Cho K, van Merrienboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, Bengio Y (2014) Learning phrase representations using rnn encoder–decoder for statistical machine translation. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics, Doha, Qatar, pp 1724–1734

7. Cornuel E (2005) A vision for business schools, vol 24. Emerald Group Publishing

8. Cortez P, Embrechts MJ (2013) Using sensitivity analysis and visualization techniques to open black box data mining models. Information Sciences 225:1–17

9. Dellermann D, Ebel P, Söllner M, Leimeister JM (2019) Hybrid intelligence. Business & Information Systems Engineering pp 1–7

10. Dernoncourt F, Lee JY (2017) Pubmed 200k rct: a dataset for sequential sentence classification in medical abstracts. In: Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers), vol 2, pp 308–313

11. Dernoncourt F, Lee JY, Szolovits P (2017) Neural networks for joint sentence classification in medical paper abstracts. In: Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, vol 2, pp 694–700

12. Di Bitetti MS, Ferreras JA (2017) Publish (in english) or perish: The effect on citation rate of using languages other than english in scientific publications. Ambio 46(1):121–127

13. Fawcett T (2006) An introduction to ROC analysis. Pattern Recognition Letters 27:861–874

14. Feinerer I, Buchta C, Geiger W, Rauch J, Mair P, Hornik K (2013) The textcat package for n-gram based text categorization in r. Journal of statistical software 52(6):1–17

15. Glorot X, Bordes A, Bengio Y (2011) Deep sparse rectifier neural networks. In: Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, pp 315–323

16. Goodfellow I, Bengio Y, Courville A, Bengio Y (2016) Deep learning, vol 1. MIT press Cambridge

17. Hossain M (2012) Users' motivation to participate in online crowd-sourcing platforms. In: Innovation Management and Technology Research (ICIMTR), 2012 International Conference on, IEEE, pp 310–315
18. Khabsa M, Giles CL (2014) The number of scholarly documents on the public web. PloS one 9(5):e93949
19. Kim Y (2014) Convolutional neural networks for sentence classification. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp 1746–1751
20. Kitchenham B, Brereton P (2013) A systematic review of systematic review process research in software engineering. Information and software technology 55(12):2049–2075
21. LeCun Y, Kavukcuoglu K, Farabet C (2010) Convolutional networks and applications in vision. Proceedings of 2010 IEEE International Symposium on Circuits and Systems pp 253–256
22. Liu Y, Wu F, Liu M, Liu B (2013) Abstract sentence classification for scientific papers based on transductive svm. Computer and Information Science 6(4):125
23. Massung E, Coyle D, Cater KF, Jay M, Preist C (2013) Using crowdsourcing to support pro-environmental community activism. In: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, ACM, pp 371–380
24. Michalska-Smith MJ, Allesina S (2017) And, not or: Quality, quantity in scientific publishing. PloS one 12(6):e0178074
25. Moro S, Laureano R, Cortez P (2011) Using data mining for bank direct marketing: An application of the crisp-dm methodology. In: Proceedings of European Simulation and Modelling Conference-ESM'2011, EUROSIS-ETI, pp 117–121
26. Moro S, Cortez P, Rita P (2015) Business intelligence in banking: A literature analysis from 2002 to 2013 using text mining and latent dirichlet allocation. Expert Systems with Applications 42(3):1314–1324
27. Moro S, Cortez P, Rita P (2017) A framework for increasing the value of predictive data-driven models by enriching problem domain characterization with novel features. Neural Computing and Applications 28(6):1515–1523
28. Moro S, Ramos P, Esmerado J, Jalali SMJ (2019) Can we trace back hotel online reviews characteristics using gamification features? International Journal of Information Management 44:88–95
29. Morschheuser B, Hamari J, Koivisto J (2016) Gamification in crowdsourcing: a review. In: System Sciences (HICSS), 2016 49th Hawaii International Conference on, IEEE, pp 4375–4384
30. Myers MD, et al. (1997) Qualitative research in information systems. Management Information Systems Quarterly 21(2):241–242
31. Ng A (2017) Machine learning yearning
32. Panettieri J (2017) Cloud market share 2017: Amazon aws, microsoft azure, ibm, google. ChannelE2E

33. Pennington J, Socher R, Manning C (2014) Glove: Global vectors for word representation. In: Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP), pp 1532–1543
34. Rebai I, BenAyed Y, Mahdi W (2016) Deep multilayer multiple kernel learning. Neural Computing and Applications 27(8):2305–2314
35. Redmon J, Divvala S, Girshick R, Farhadi A (2016) You only look once: Unified, real-time object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 779–788
36. Ribeiro MT, Singh S, Guestrin C (2016) "why should I trust you?": Explaining the predictions of any classifier. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016, pp 1135–1144
37. Seyler D, Li L, Zhai C (2018) Identifying compromised accounts on social media using statistical text analysis. arXiv preprint arXiv:180407247
38. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R (2014) Dropout: a simple way to prevent neural networks from overfitting. The Journal of Machine Learning Research 15(1):1929–1958
39. Stewart O, Lubensky D, Huerta JM (2010) Crowdsourcing participation inequality: a scout model for the enterprise domain. In: Proceedings of the ACM SIGKDD Workshop on Human Computation, ACM, pp 30–33
40. Stoean R (2018) Analysis on the potential of an ea–surrogate modelling tandem for deep learning parametrization: an example for cancer classification from medical images. Neural Computing and Applications pp 1–10
41. Suciu G, Scheianu A, Vochin M (2017) Disaster early warning using time-critical iot on elastic cloud workbench. In: Black Sea Conference on Communications and Networking (BlackSeaCom), 2017 IEEE International, IEEE, pp 1–5
42. Torres JP, de Piñerez Reyes RG, Bucheli VA (2018) Support vector machines for semantic relation extraction in spanish language. In: Colombian Conference on Computing, Springer, pp 326–337
43. Tsapatsoulis N, Djouvas C (2019) Opinion mining from social media short texts: Does collective intelligence beat deep learning? Front Robotics and AI DOI 10.3389/frobt.2018.00138, URL https://doi.org/10.3389/frobt.2018.00138
44. Welinder P, Perona P (2010) Online crowdsourcing: rating annotators and obtaining cost-effective labels. In: Computer Vision and Pattern Recognition Workshops (CVPRW), 2010 IEEE Computer Society Conference on, IEEE, pp 25–32
45. Witten I, Frank E, Hall M, Pal C (2017) Data Mining: Practical Machine Learning Tools and Techniques, 4th edn. Morgan Kaufmann, San Franscico, USA, San Francisco, CA
46. Xia Y, Liao X, Li S, Shi Q, Wu J, Li K (2019) Extracting PICO elements from RCT abstracts using 1-2gram analysis and multitask classification. CoRR abs/1901.08351, URL http://arxiv.org/abs/1901.08351, 1901.08351
47. Zhang Q, Yang LT, Chen Z, Li P (2018) A survey on deep learning for big data. Information Fusion 42:146–157

48. Zhang X, Zhao J, LeCun Y (2015) Character-level convolutional networks for text classification. In: Advances in neural information processing systems, pp 649–657

49. Zheng H, Li D, Hou W (2011) Task design, motivation, and participation in crowdsourcing contests. International Journal of Electronic Commerce 15(4):57–88

50. Zubiaga A, Kochkina E, Liakata M, Procter R, Lukasik M, Bontcheva K, Cohn T, Augenstein I (2018) Discourse-aware rumour stance classification in social media using sequential classifiers. Inf Process Manage 54(2):273–290, DOI 10.1016/j.ipm.2017.11.009, URL https://doi.org/10.1016/j.ipm.2017.11.009