# Argument Mining

*Author:*
Klio Fragkedaki

*Supervisor:*
Prof. Panagiotis Louridas

*An assignment submitted as part of*
*Bachelor degree*

*in the*

Department of Management Science and Technology

September 30, 2019

# Contents

# Acknowledgments

I would like to thank Vasiliki Efstathiou for being the second annotator and a great supervisor, as well as my Professor Panos Louridas.

# List of Abbreviations

| | |
|---|---|
| **SVM** | **S**upport **V**ector **M**achine |
| **LR** | **L**ogistic **R**egression |
| **NB** | **N**aive **B**ayes classifier |
| **RF** | **R**andom **F**orest |
| **RNN** | **R**ecurrent **N**eural **N**etworks for Language Models |
| **RF** | **R**andom **F**orest |
| **CRF** | **C**onditional **R**andom **F**orest |
| **ML** | **M**aximum **L**ikelihood |
| **TES** | **T**extual **E**ntailment **S**uites |
| **P** | **P**arsing Using a Context-Free Grammar |
| **LSTM** | **L**ong **S**hort-**t**erm **M**emory |
| **POS** | **P**art **O**f **S**peech |

# Chapter 1

# Introduction

## 1.1 Definition

Argument mining is a relatively new research field in natural language processing. The aim of this research is the auto detection and identification of argumentative structures expressed in text. In order to perform extraction and evaluation of arguments, computer science and artificial intelligence is used.

An argument is a group of premises conducted to support a claim (Palau and Moens, 2009). When it comes to real world, arguments are hardly identified even by experts (Lippi and Torroni, 2015). The ambiguity of natural language, the implicit content, the different ways of expressing and the complex structure of arguments are the main reasons why argument mining is a challenging research field. Labeled corpora are scarce which is a fact that slows down field's potential growth (Lippi and Torroni, 2015).

The purpose of argument mining is to understand what kind of views have been expressed in the examined text and why they are held. Argument mining has derived from opinion mining and sentiment analysis research area, in which the only goal is to understand the opinions about a certain topic (Lawrence and Reed, 2015).

## 1.2 Research Goal

My research goal is to identify argumentative statements by using two different approaches; the structural approach which is based on hand coded rules and the statistical approach, which based on supervised and deep learning algorithms.

The structural approach uses lexical cues that have been identified by linguists as signs of argumentative speech. As an example, words such as "because", "therefore", "in order to" are common cues of arguments. However, these argumentative patterns are rarely used in practice, since human discourse involves a lot of information which is being implied rather than being explicitly stated.

On the other hand, the statistical approach relies on examples of pieces of text that have been manually labeled as argumentative or non-argumentative. These are used for training models in order to automatically identify arguments in free text without the use of predefined lexical cues and rules. The challenging part is the construction of a manually annotated data-set, given the fact that a large amount of data are required for training such models.

The fundamental research questions that will be addressed in this assignment are the following:

- To what extent are the lexical rules drafted by a structural approach capable of successfully identifying arguments in existing resources of labeled data?

- Do the statistical approaches outperform these results?

## 1.3    Assignment's Structure

This paper of research is organized into 6 chapters. Chapter 2 presents the state of the art in argument mining, and introduces the two different approaches; the structural and the machine learning approach. Chapters 3 and 4 describe in detail the methods and results of both approaches implemented in the scope of this study. Chapter 5 contains the corpora created for the supervised algorithm, while chapter 6 concludes with a look to future work.

# Chapter 2

# State-of-the-Art

Arguments do not have a universally accepted definition; though there are plenty of well-described proposals. According to (Walton, 2009), an argument is a group of statements which splits into three portions, which are conclusion, set of premises, and an inference leading from premises to conclusion. These concepts have been widely accepted in literature, but they are defined in slightly different ways. Conclusions are also referred to as claims, premises as evidence or reasons, while the link between claims and evidence is the argument (Lippi and Torroni, 2015).

A claim is supported or argued by one or more premises and it is the main part of an argumentative text. Claims are controversial in terms of validity and need premises to endorse readers' acceptance (Stab and Gurevych, 2014). Argumentation schemes and their common patterns provide a way to both identify and determine arguments (Lawrence and Reed, 2015).

The term of argumentation used to be connected with the process of argument construction (Lippi and Torroni, 2016). After the emergence of text mining procedures, this term defines the process of argument identification in text (Lippi and Torroni, 2016). The research field of argument mining is about the automatic recognition of argumentative structures expressed in natural language texts. Argument mining utilizes methods and techniques used in natural language processing, such as machine learning and sentiment analysis (Lippi and Torroni, 2015).

In general, argument mining procedure is separated into linguistic and computational part, as described in figure 2.1. Regarding the linguistic part, large corpora of manually annotated argument data are being created based on a common agreement among annotators about argument's structure. On the other hand, the computational part is separated into two main styles of automation, the structural and the statistical approach. (Budzynska and Villata, 2015)

In **structural or grammar approach**, linguists aim to retrieve lexical patterns, rules or categories while annotating a training corpus. For example, it might be noticed that words like "because", "since", "however" are signs of arguments inside a specific corpus (Budzynska and Villata, 2015). These signs are called indicators, and point out the connection between claims and premises inside a text (Lawrence and Reed, 2015). Indicators are declared as linguistic expressions that connect statements and provide an unambiguous recognition of argumentative structure (Webber, Egg, and Kordoni, 2012).

A lot of research has been applied in order to be found words and expressions revealing argumentative structure (Van Eemeren, Houtlosser, and Henkemans, 2007, Knott and Dale, 1994). Apart from indicators, other structural techniques have been applied for argument mining. Such techniques are argumentation schemes (Feng and Hirst, 2011), dialogical context (Budzynska et al., 2014), and semantic context (Cabrio and Villata, 2012) or a combination of them.
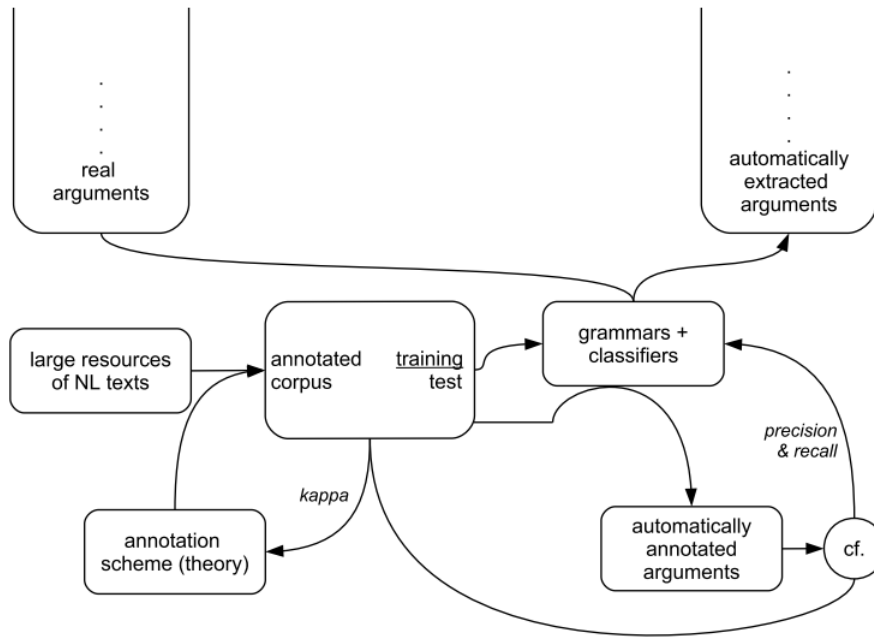
FIGURE 2.1: Natural language processing techniques
**Source:** Budzynska and Villata, 2015

In **statistical approach**, linguists are replaced by algorithms. These algorithms are basically classifiers developed for automating the argument annotation procedure (Budzynska and Villata, 2015). The first attempts for the before mentioned automation were made in (Moens et al., 2007), in which text is separated into sentences, and then each sentence is classified as argumentative or non-argumentative based on its lexical or syntactic features. As a result, (Palau and Moens, 2009) presented an additional separation of argumentative sentences as premises or conclusions. As regards the automatic recognition of argumentative schemes, it was introduced in (Walton, 2011) and it was based on the idea of connecting each scheme with a group of indicators. The paper's proposal is first indicating the arguments included in text, and then matching them to a given list of argument schemes. (Feng and Hirst, 2011) classifies annotated argumentation structures into a list of five common argumentation schemes. In (Lippi and Torroni, 2015), the authors describe a framework for claim detection in unstructured data-sets without any contextual information. Because arguments are often expressed through rhetorical structures, the previously mentioned framework was built based on an SVM classifier which captures similarities among parse trees via Tree Kernels. This method is used for measuring likeliness of two trees regarding their common substructures. Furthermore, Habernal and Gurevych (Habernalt and Gurevych, 2016) try to evaluate argument convincingness by assessing their qualitative properties. Using an annotated corpus of 26,000 sentences, their purpose is to predict which argument is more convincing between a pair of arguments and to rank arguments regarding the topic and their convincingness, through the usage of SVM and LSTM algorithms.

Various traditional machine learning algorithms have been employed in the context of argument mining (Figure 2.2). More specifically, most of the algorithms that have been implemented are Support Vector Machines (Mochales and Moens, 2011; Park and Cardie,

2015; Stab and Gurevych, 2014; Eckle-Kohler, Kluge, and Gurevych, 2015), Logistic Regression (Levy et al., 2014; Rinott et al., 2015), Naive Bayes classifiers (Mochales and Moens, 2011; Biran and Rambow, 2011; Park and Cardie, 2015; Eckle-Kohler, Kluge, and Gurevych, 2015), Maximum Entropy classifiers (Mochales and Moens, 2011), and Decision Trees and Random Forests (Stab and Gurevych, 2014; Eckle-Kohler, Kluge, and Gurevych, 2015). All mentioned classifiers are trained in labeled corpora. Thus, some parts of the annotated text are given, alongside with the associated label, and during training stage a model is being produced. This model is used to perform predictions on new unlabeled text. (Lippi and Torroni, 2016)

| System | SC | | | | | | | BD | | SP | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SVM | LR | NB | ME | DT | RF | RNN | CRF | ML | TES | P | SVM | NB |
| Eckle-Kohler et al. [2015] | X | | X | | | X | | | | | | | |
| Lippi and Torroni [2015] | X | | | | | | | | | | | | |
| Rinott et al. [2015] | | X | | | | | | | | | | | |
| Sardianos et al. [2015] | X | | | | | | X | X | | | | | |
| Boltuzic and Snajder [2014] | | | | | | | | | | X | | X | |
| Goudas et al. [2014] | X | | | | | | | X | | | | | |
| Levy et al. [2014] | | X | | | | | | | X | | | | |
| Stab and Gurevych [2014b] | X | | X | | X | X | | | | | | X | |
| Cabrio and Villata [2012a] | | | | | | | | | | X | | | |
| Rooney et al. [2012] | X | | | | | | | | | | | | |
| Biran and Rambow [2011] | | | X | | | | | | | | X | | X |
| Mochales Palau and Moens [2011] | X | | X | X | | | | | | | X | | |

FIGURE 2.2: Machine learning algorithms that have been used for argument mining
**Source:** Lippi and Torroni, 2016

Despite the fact that researchers have tried to make a comparison between these algorithms, there is no clear proof of which classifier is more appropriate for argumentation mining. In fact, most of the research efforts have been settled down on finding appropriate features for improving performance instead of implementing new specifically designed models and algorithms for solving argument identification problem (Lippi and Torroni, 2016).

To sum up, a number of different approaches have been applied to argument identification problem. The research community solutions are ranging from linguistic techniques (Garcia Villalba and Saint-Dizier, 2012) and topic modeling (John Lawrence, Chris Reed, Colin Allen, Simon McAlister, Andrew Ravenscroft, 2014), to supervised machine learning algorithms( firstly implemented by Moens et al., 2007).

# Chapter 3

# Methods

In this research paper, we attempt to apply two different approaches for recognizing argumentative sentences. These approaches cover both a structured methodology, which is related to the selection of hand-coded linguistic rules, and a statistical one, that includes the implementation of supervised algorithms; namely, Random Forest classifier and sequence classification with LSTM.

## 3.1  Structural Approach

The structural approach is based on lexical cues, rules or patterns for identifying arguments inside a given text. These cues are also referred to as argument indicators, since they are connecting claims and premises, signaling argumentative relations.

| Argumentative Indicators based on (Knott and Dale, 1994) | | | | | | | |
|---|---|---|---|---|---|---|---|
| Indicator | POS | Indicator | POS | Indicator | POS | Indicator | POS |
| even though | none | first | adv | against | none | last | adv |
| naturally | none | most | {"[a-z]*ly": "adv"} | if | none | (T\|t)(he more).+?(the more) | none |
| once more | none | more | {"[a-z]*ly": "adv"} | once again | none | (T\|t)(he more).+?(the less) | none |
| surely | none | second | adv | so | mark | third | adv |
| should say | none | too | (too)($\|[\\.]) | might say | none | may say | none |
| could say | none | while | mark | as a start | none | in order to | none |
| still | adv | that is | none | since | mark | yet | (Y\|y)(et)[ˆ\.]. |
| that | mark | above all | none | actually | none | after all | none |
| afterwards | none | all in all | none | also | none | although | none |
| anyway | none | as a consequence | none | as a result | none | at any rate | none |
| at first blush | none | at first view | none | at the outset | none | because | none |
| by comparison | none | by the same token | none | certainly | none | consequently | none |
| correspondingly | none | despite the fact that | none | either | none | equally | none |
| even then | none | every time | none | except insofar as | none | firstly | none |
| for a start | none | for instance | none | further | none | for the simple reason | none |
| accordingly | none | admittedly | none | after that | none | all the same | none |
| alternatively | none | always assuming that | none | as | none | as a corollary | none |
| at first | none | at first sight | none | at the moment when | none | at the same time | none |
| but | none | by contrast | none | by the way | none | clearly | none |
| conversely | none | despite that | none | essentially | none | even so | none |
| eventually | none | except | none | finally | none | first of all | none |
| for example | none | for one thing | none | for this reason | none | furthermore | none |
| hence | none | in actual fact | none | in any case | none | in conclusion | none |
| in fact | none | in other words | none | in short | none | in sum | none |
| incidentally | none | instead | none | merely because | none | just as | none |
| meanwhile | none | it might appear that | none | as long as | none | as well | none |
| notably | none | moreover | none | of course | none | nevertheless | none |
| on one hand | none | not only | none | now that | none | no doubt | none |
| on the grounds that | none | on the assumption that | none | on the one side | none | on the other side | none |
| plainly | none | otherwise | none | so that | none | providing that | none |
| such that | none | secondly | none | sure enough | none | simply because | none |
| thereafter | none | summing up | none | therefore | none | suppose that | none |
| thirdly | none | the fact is that | none | to be sure | none | though | none |
| to sum up | none | to conclude | none | undoubtedly | none | to take an example | none |
| whenever | none | to the extent that | none | whereas | none | what is more | none |
| wherever | none | for the reason that | none | besides | none | (E\|e)(ither).+?(or) | none |
| in one hand | none | (N\|n)(either).+?(nor) | none | on one side | none | in this case | none |
| in point of fact | none | as a matter of fact | non | provided that | none | presumably | none |
| rather than | none | regardless | none | as an example | none | simply | none |
| in order that | none | | | | | | |

A list of indicators were extracted from the corpus created by (Knott and Dale, 1994). This corpus includes often-used words or phrases in arguments according to paper's authors. Based on these words, a dictionary was developed containing as keys the extracted words, and as values, their specific part of speech in argumentative sentences. It needs to

be mentioned that words, considered by us as usual or non-usual in argumentative structures, were added or removed respectively from the dictionary. For this purpose, there were created five methods in Python for paper's extraction, modification, as well as dictionary's creation (Appendix 6). The indicators that demonstrate the previously referred dictionary is presented in the table above.

Apart from dictionary's development, a way to handle and encapsulate corpora into the same format was necessary, and the code developed for this purpose is shown in Appendix 7. Each data-set was differently displayed, from unstructured text to sentence labeled data. This is the reason why there was created a *datasets.ini* file containing information about data, for example the number of column indicating the sentence or/and the label, which sheet includes the desired data, or which is the data-set's path. The key of each record was the name of every corpora as it was saved in local file. So, depending on the data-set's type (excel, csv or txt file) and its configurations, other actions were applied in order to returned a list of sentences and their labels in case corpora was annotated.

By using the previously created dictionary and corpora handler, argument identification had to take place (Appendix 8). For this reason, part of speech tagging was necessary, so as a sentence's words and their POS to be compared to those words included in the dictionary. Statements tokenization was achieved through the usage of a library called *spaCy*, which is an open-source NLP library written in Python and Cython, and it was selected due to its performance and efficiency comparing to other libraries. If any of matches between the dictionary and a given sentence occur, the sentence is characterized as argumentative, otherwise as non-argumentative. As regards the labeled corpora, the algorithm's outcomes and the given labels, which is considered to be the truth, are correlated so as four counters to be calculated; False Positives, False Negatives, True Positives and True Negatives. These counters are used for measuring accuracy that is a metric for reviewing algorithm's results, which will be used and described at Chapter 5.

## 3.2 Statistical Approach

The statistical approach is relying on algorithms developed for automating the argument annotation process. These algorithms receive as input, a representation of data that can be understood by computers. For this reason, textual data of our corpora need to be transformed into numeric tensors. (Chollet, 2017)

Text is considered as a form of sequence data, and tokens are the different units into which a text can be split. These units can be either words, characters or n-grams (Chollet, 2017). After text tokenization, the next step is connecting numeric vectors with the occurred tokens. Then, the sequences of vectors instead of words is fed into the selected algorithm.

There are a variety of ways for vector and token association. Two and the most known ones for forms of sequence data are the one-hot encodingText and word-embedding methods (Chollet, 2017), which will be used in the implemented algorithms later on.

As regards one-hot encodingText is a basic way to transform tokens into vectors. Every word is connected with a unique integer index, which is turned into a binary vector of vocabulary's size. Another version of this method is the one-hot hashing trick, which is used when vocabulary's unique tokens are too much that is difficult to be handled. This approach hashes words into a fixed-sized vector through a hashing function, rather than allocating an index to each word, and then creating a dictionary for these indexes. The disadvantages of this method is the hash collisions, and in general the inability of word's correlation that leads to high dimentionality. (Chollet, 2017)

Word embedding, or dense word vectors, is another popular way to associate vectors with words. In contrast with one-hot encodingText, this technique has low-dimensionality

and float vectors. Furthermore, dense word vectors are trained by the data provided, and thus packing more and more data into the same dimensions. This is achieved by minimizing the geometric distance between related words, like synonyms, while maximizing it when words have different semantic (Chollet, 2017).

The implemented algorithms for argument mining problem are the Random Forest classification algorithm, using one-hot encodingText, and the LSTM-RNN algorithm, using word embedding method. The annotation of argumentative and non-argumentative sentences can be considered as a binary classification in Random Forest (Stab and Gurevych, 2014) or a statistical structure mapping of written language in LSTM (Chollet, 2017).

### 3.2.1 Random Forest Classification Algorithm

The Random Forest algorithm is a collection of decision trees, where each one of them are slightly different from the others either by selecting the data-points or by selecting different features (Müller and Guido, 2016). Every decision tree, built during the training period, predicts to which class a specific input is more suitable (argument or non-argument). After this process is completed, the chosen prediction for each sentence is the one that has the majority of votes across the decision trees. Chollet, 2017

By using the `sklearn.ensemble` library, a Random Forest classification algorithm was implemented (Appendix 9). First of all, an annotated corpora containing both argumentative and non-argumentative sentences is loaded and features regarding how many words, punctuation characters and uppercase characters were created (see Chapter 5). Afterwords, the sentences were tokenized in order to be transformed into a format that computer can understand, as mentioned before. For the tokenization `OneHotEncoder` of `sklearn.preprocessing` library was used. The corpora was seperated into a training and a test set of data, and then model was defined and trained using the train data. Finally, the model is tested and a matrix showing the real and the predicted data was presented in a heatmap plot.

The model is applied by using two different methods for training. The one method is using the features created, like in (Lawrence and Reed, 2016), while the other is by using the tokenized sentences.

### 3.2.2 LSTM-RNN Algorithm

Deep learning in natural-language processing problems is aiming to identify patterns of words or sentences (Chollet, 2017). The most appropriate deep learning model for sequence data processing, like sequences of words in our case, is the Recurrent Neural Network algorithm. RNN is a neural network algorithm that is implementing internal iterations. The RNN resets its state every time a new independent process occurs, which means that each sequence is a different data input to the network. However, network include loops over each sequence's elements. RNN includes the so-called LSTM layer. The Long Short-Term Memory algorithm was developed by Hochreiter and Schmidhuber in 1997, and the purpose of creation was to solve the "vanishing-gradient problem" of SimpleRNN algorithm. LSTM is basically allowing previous information to be re-injected across different timesteps. (Chollet, 2017)

By using the tensorFlow keras API, an LSTM algorithm was conducted and implemented (Appendix 10). First of all, the corpora created in Chapter 4 is loaded, and data applied to one list of sentences and one list of labels. Following, the corpora is split into a training and a testing data-set. Furthermore, the lists of trained sentences were tokenized, taking into account only the first ten-thousand most frequent words. The lists of the occurred trained sequences, which are basically lists of integers after the tokenization process, are transformed into a two-dimensions numpy array that has a shape of (number of

sequences, number of timesteps). The number of timesteps is basically the length of the longest sequence. It has to be mentioned that shorter sequences are padded with values at the end, so every sequence has the same shape.

As regards word embedding, Global Vectors for Word Representation (*Glove*) was used, which is a file containing 100-dimensional pre-computed embedding vectors for 400,000 of English words. This file helped in building an index that associates words with number vectors, and through which an embedding matrix of shape (max_words= 10000, embedding_dim = 100) was created and loaded in Embedding layer later on. The Embedding layer is a dictionary-like that maps integer indexes for certain words to dense vectors, and it is the first layer of our model. The second layer applies to LSTM method, while the third layer of Dense was added to end our stack of layers with an equal number of units and classes created. After constructing our model, we are freezing the Embedding layer to avoid deleting what has already been learned.

Finally, the model was trained based on the training data, which are explicitly separated into training and validation samples used for learning word embedding based on our corpora. The last step is testing the model on the test data by firstly tokenizing the sentences, and then evaluating the results occurred.

# Chapter 4

# Data Curation

In order to successfully apply the statistical learning approach, a well-structured training data-set is needed. In this section, the process of data curation is elaborated so as both argumentative and non-argumentative sentences to be found.

Most of the argumentative sentences included in our corpora were found on two IBM data-sets created for this purpose (Table 4.1), and are about three main topics; Video Games, Democracy and Multiculturalism. It has to be mentioned that some of the data were duplicated, and thus Python code of Appendix 2 was created so as to be removed.

| Source Data | file | Total Data | Duplicated | Arguments | Non-Arguments |
|---|---|---|---|---|---|
| Aharoni et al., 2014a | CDEdata.xls | 1292 | 967 | 325 | - |
| Bar-Haim et al., 2017 | claim_stance_dataset_v1 | 2394 | 56 | 2366 | - |

TABLE 4.1: Argumentative Data Used

However, the exploitation of the existing annotated data-sets regarding argument detection has the obstacle of lacking non-argumentative instances. The previously mentioned IBM corpora contain only phrases that have been manually annotated as positive instances of arguments. This means that it is impossible to train a supervised algorithm in classifying non-arguments without any negative examples.

So, our purpose was to gather an equal number of argumentative and non-argumentative sentences that have the same context with the curated data of IBM. Therefore, plain text referring to Video Games was found in additional IBM data-sets. These raw data files were split into sentences, and each of these sentences was labeled as argument or non-argument by authors of this research paper, and the Table 4.2 depicts the results. The data referred as "Not used" were blank or incomplete lines.

The non-arguments collected were not enough, thus it was decided to scrap data from Wikipedia articles. The topics of these articles were similar to the previously gathered data, and more specifically about Video games, Democracy and Multiculturalism. The code used for the scraping process, as weel as the scraped pages, are aligned at Appendix 1.

| Source Data | file | Total Data | Not Used | Arguments | Non-Arguments |
|---|---|---|---|---|---|
| Mirkin et al., 2017 | asr/DJ_1_ban-video-games_pro.wav.asr.txt | 9 | 3 | 5 | 1 |
| Mirkin et al., 2017 | asr/EH_1_ban-video-games_pro.wav.asr.txt | 20 | 3 | 12 | 5 |
| Mirkin et al., 2017 | asr/HE_1_ban-video-games_pro.wav.asr.txt | 21 | 1 | 12 | 8 |
| Mirkin et al., 2017 | asr/SN_1_video-games_pro.wav.asr.txt | 28 | 10 | 15 | 3 |
| Mirkin et al., 2017 | asr/TL_1_ban-video-games_pro.wav.asr.txt | 19 | 6 | 8 | 5 |
| Mirkin et al., 2017 | asr/YB_1_ban-video-games_pro.wav.asr.txt | 19 | 4 | 7 | 8 |
| Aharoni et al., 2014a | wiki12_articles/Gender_representa-tion_in_video_games | 39 | 4 | 5 | 30 |

TABLE 4.2: Argumentative and Non-Argumentative Data Used

Assuming that Wikipedia articles' authors are objective and do not express their point of view, most of the data scraped were included as non-arguments in the data-set. It has to be mentioned that the data collected from this procedure were previously checked from the python code-described in the Chapter 3 ( Appendix 6). The sentences classified as non-arguments were added to the created data-set, while the others were considered as ambiguous (Table 4.3).

| Topic of Wikipedia | Total Data | Not Used | Controversial sentences | Non-Arguments |
|---|---|---|---|---|
| Early Hstory of video games | 143 | 19 | 59 | 65 |
| Fourth generation of video game consoles | 65 | 6 | 32 | 27 |
| Game Boy | 84 | 22 | 23 | 39 |
| Game design | 223 | 32 | 71 | 120 |
| Game | 191 | 28 | 89 | 74 |
| Gaming Computer ' | 129 | 8 | 62 | 59 |
| Gaming disorder | 12 | 6 | - | 6 |
| History of video games | 604 | 70 | 224 | 310 |
| Home computer | 359 | 236 | 54 | 69 |
| Nintendo | 393 | 70 | 133 | 190 |
| PC game | 248 | 60 | 87 | 101 |
| Video game | 434 | 82 | 154 | 198 |
| Video game addiction in China | 58 | 4 | 9 | 45 |
| Video game addiction | 257 | 138 | 70 | 49 |
| Video game console | 337 | 261 | 41 | 35 |
| Video game culture | 292 | 74 | 104 | 114 |
| Video game development | 446 | 229 | 102 | 115 |
| Video game industry | 275 | 49 | 91 | 135 |
| Video game music | 460 | 176 | 146 | 138 |
| Video game programmer | 164 | 19 | 72 | 82 |
| Video game-related health problems | 52 | 7 | 22 | 23 |
| Video gaming in Japan | 300 | 102 | 82 | 116 |
| Video gaming in the United States | 119 | 31 | 27 | 61 |
| The Game Awards | 36 | 2 | 18 | 16 |
| Multicultural transruption | 45 | 3 | 20 | 22 |
| Multicultural and diversity management | 41 | 7 | 17 | 17 |
| Multicultural education | 248 | 28 | 104 | 116 |
| Multiculturalism in Australia | 156 | 37 | 55 | 54 |
| Criticism of multiculturalis | 237 | 56 | 96 | 85 |
| Cultural pluralism | 28 | 7 | 12 | 9 |
| Multiculturalism in Canada | 205 | 28 | 84 | 93 |
| Multiculturalism | 449 | 62 | 160 | 227 |
| Democracy Index | 59 | 18 | 17 | 24 |
| Direct democracy | 162 | 22 | 59 | 81 |
| Types of democracy | 25 | 7 | 9 | 9 |
| Representative democracy | 56 | 7 | 26 | 23 |
| Criticism of democracy | 191 | 102 | 55 | 34 |
| Athenian democracy | 335 | 41 | 147 | 147 |
| History of democracy | 394 | 81 | 126 | 187 |
| Democracy | 452 | 71 | 193 | 188 |

TABLE 4.3: Non-Argumentative Data Used & Controversial Sentences

The previously described data were concatenated (Appendix 3), so as to be used in the statistical approach algorithms. The corpora is composed by both arguments and non-arguments (Table 4.3).

| Argumentative and No-Argumentative Data Used | | | | |
|---|---|---|---|---|
| File | Total Data | Arguments | Non-Arguments | False-Positive Arguments |
| dataset.csv | 6318 | 2755 | 3563 | - |
| found_fp.csv | 2952 | - | - | 2952 |

TABLE 4.4: Data included in our corpora

As regards the ambiguous sentences of Wikipedia articles that was mentioned before, they were all saved in a file named `found_ambiguous.csv` (Appendix 4). This file indicates all the keywords responsible for these controversial results. The indicators- described in the previous Chapter and found in Wikipedia articles- were counted using code of Appendix 5, and the following table was the outcome of that enumeration.

| Indicator Found | Number of Sentences | Indicator Found | Number of Sentences |
|---|---|---|---|
| **as** | **968** | **that** | **542** |
| **also** | **440** | **because** | **121** |
| **as well** | **101** | **for example** | **92** |
| (E\|e)(ither).+?(or) | 44 | as a result | 23 |
| **because** | **121** | notably | 13 |
| for instance | 17 | **but** | **326** |
| that is | 40 | **while** | **197** |
| actually | 26 | against | 3 |
| still | 85 | so that | 12 |
| though | 87 | besides | 4 |
| furthermore | 13 | eventually | 41 |
| **more +** | **40** | either | 53 |
| clearly | 5 | **if** | **107** |
| since | 34 | on the grounds that | 3 |
| although | 85 | third | 6 |
| consequently | 4 | as a consequence | 2 |
| rather than | 66 | instead | 42 |
| nevertheless | 9 | except | 5 |
| otherwise | 12 | in fact | 10 |
| too | 3 | not only | 22 |
| on one side | 1 | simply | 26 |
| moreover | 9 | hence | 4 |
| therefore | 25 | every time | 1 |
| in other words | 4 | despite the fact that | 2 |
| in order to | 43 | as | 9 |
| at the same time | 7 | of course | 2 |
| finally | 13 | as an example | 2 |
| first | 29 | even though | 14 |
| lastly | 3 | equally | 7 |
| whereas | 13 | (T\|t)(he more).+?(the more) | 2 |
| regardless | 7 | for this reason | 2 |
| simply because | 1 | by contrast | 3 |
| naturally | 4 | at any rate | 1 |
| in short | 2 | in this case | 3 |
| such that | 5 | essentially | 5 |
| (N\|n)(either).+?(nor) | 4 | whenever | 4 |
| second | 5 | as long as | 2 |
| even then | 3 | as a matter of fact | 1 |
| accordingly | 3 | provided that | 1 |
| conversely | 3 | alternatively | 3 |
| afterwards | 6 | thereafter | 1 |
| meanwhile | 10 | once again | 3 |
| once more | 1 | above all | 1 |
| by comparison | 1 | surely | 2 |
| undoubtedly | 2 | on the one side | 1 |
| at first | 1 | presumably | 2 |
| after all | 1 | what is more | 1 |
| certainly | 1 | anyway | 1 |
| so | 16 | most | 22 |
| **further** | **53** | | |

TABLE 4.5: Indicators found in Wikipedia articles

The goal of this process was to test in Wikipedia articles the argumentative indicators included in the Chapter 3's dictionary, and upon cross-examination to remove indicators that do not usually reveal argumentative sentences. Based on the Table 4.5, the indicators marked as bold are the ones found in the majority of Wikipedia sentences. More than 10% of each indicator's sentences was examined by the authors of this research papers, and the keywords, that did not usually pointed out argumentative statements and removed from the dictionary, are represented in the following. The annotations are described in detail inside the file `Results/found_ambiguous_results_annotated.xlsx`.

| Indicator | Total Number of Sentences | Number of Sentences Examined | Arguments | Non-Argument | Depending on the context | Removed from the dictionary |
|---|---|---|---|---|---|---|
| as | 968 | 131 | 22 | 98 | 11 | Yes |
| that | 542 | 91 | 20 | 54 | 17 | Yes |
| also | 440 | 64 | 8 | 46 | 10 | Yes |
| but | 326 | 50 | 7 | 36 | 7 | Yes |
| while | 197 | 35 | 6 | 24 | 5 | Yes |
| because | 121 | 33 | 20 | 10 | 3 | No |
| if | 107 | 20 | 6 | 12 | 2 | No |
| as well | 101 | 20 | 5 | 14 | 1 | No |
| for example | 92 | 8 | 5 | 3 | 0 | No |
| further | 53 | 20 | 4 | 15 | 1 | Yes |
| against | 51 | 21 | 3 | 16 | 2 | Yes |
| more +a dverb in 'ly' | 40 | 16 | 4 | 12 | 0 | No |

TABLE 4.6: Indicators found in Wikipedia articles

It was also observed that the following indicators were used in combination with punctuation, and that's why they were modified into the formats of the table 4.7.

| Indicator | Modified to |
|---|---|
| while | , while |
| that is | , that is, |
| still | , still |
| as a result | as a result, |

TABLE 4.7: Indicators found in Wikipedia articles

# Chapter 5

# Results

Results for both structured and statistical implementations presented in Chapter 3, are applied to a set of corpora in order to be evaluated.

## 5.1 Structural Approach

The structural approach described in previous chapters is being assessed into this section. For this purpose, the code of Appendix 8 alongside with two IBM corpora was applied. It has to be mentioned that these two IBM corpora include only argumentative sentences, and we aim to check how many of those arguments will be identified correctly by our algorithm.

The metric of accuracy was used in order to evaluate the indicators selected for recognizing argumentative sentences, while precision, recall and f1 score did not have value because of the non-arguments lack. For measuring accuracy, four counters were used; **true positives** (tp) is counting the times both algorithm and analyst labeled a sentence as argumentative, **true negatives** (tn) how often both algorithm and analyst labeled as non-argumentative, **false positive** (fp) the times the algorithm assigned as argumentative a sentence that expert recognized as non-argumentative, **false negative** (fn) how many times human identified a sentence as argumentative while algorithm did not.

- **Accuracy** represents the percentage of correctly classified sentences:

$$A = \frac{tp + tn}{tp + tn + fp + fn}$$

- **Precision** indicates the times of correctly identification instances:

$$P = \frac{tp}{tp + fp}$$

- **Recall** measures the times algorithm missed out arguments:

$$R = \frac{tp}{tp + fn}$$

- **F1 Score** presents the mean of precision and recall:

$$F = \frac{2 * P * R}{P + R}$$

| Source Data | file | Accuracy |
|---|---|---|
| Bar-Haim et al., 2017 | claim_stance_dataset_v1 | 17.50% |
| Aharoni et al., 2014b | CDEdata.xls1.00.1480.25 | 14.81% |

TABLE 5.1: Results of Structural Approach

Structural approach seems not to be able to capture a variety of argumentative structures. This is beacause, argumentative patterns are rarely used in practice, since human discourse involves a lot of information which is being implied rather than being explicitly stated.

## 5.2 Statistical Approach

### 5.2.1 Random Forest Algorithm

Supervised machine learning algorithms need a number of labeled data in order to be trained. That was the reason corpora of Chapter 4 was created, and used in the implantation of Random Forest classifier algorithm (Appendix 9). A number of 33% of the data used to train the model, while the rest of them to evaluate the results.

It has to be mentioned that two methods were used for argument's classification. The first one was by tokenizing the sentences, so as to be in a format that a computer can understand, and then training the model based on the tokenized sentences. This method had an accuracy of **56.83%**, precision **100%**, recall **0.55%**, f1 score **71.13%**, and the results are displayed in the heatmap of figure 5.1 (A). Based on the results, it seems that Random Forest trained by tokenized sentencses is not recognizing argumentative sentences. That's why we implemented another technique in which we determined some features of each sentence ,and then feed the algorithm with these features instead of the sentences. In this way the accuracy increased to **79.42%**, while precision is **74.08%**, recall **78.72%**, f1 score **76.33%** , and it deprecates to the heatmap of figure 5.1 (B). The results of the first approach are not that high, and that is probably because of the unique words.

The features used for classification are the following based on the paper (Lawrence and Reed, 2016):

- **Word Counter**: the number of words in a sentence

- **Uppercase Characters Counter**: the number of uppercase characters found

- **Punctuation or Special Characters Counter**: the number of presence punctuation characters like " "



(A) Classification based on the sentences

(B) Classification based on features

FIGURE 5.1: Predicted and Real results of Random Forest classification algorithm

As the results of table 5.1 reveal, Random Forest implements well when features are used for each of sentence. Machine learning approach seems to be a better fit for identifying arguments than lexical rules.

### 5.2.2 LSTM-RNN Algorithm

The LSTM-RNN algorithm described in previous chapters is being assessed into this section. For this purpose, the code of Appendix 10 alongside with the data-set described in Chapter 4 was executed.

The data-set used contains an equal number of argumentative and non-argumentative sentences, as well as their labels. The twenty percent of the data were used for training the model, while the rest of them for evaluating it. The model's performance over time is represented in the following plots by using the metrics of accuracy and loss.



(A)                                    (B)

FIGURE 5.2: Training and validation accuracy and loss when using pretrained word embeddings

After testing the algorithm in test data, the following results occurred with an accuracy of **85.36%**, precision **78.15%**, recall **88.90%**, f1 score **82.61%**. These results lead to a conclusion that LSTM-RNN algorithm seems to be more appropriate for argument mining problems, comparing to Random Forest and the Structural approach examined in the previous section.

```
Found 10823 unique tokens.
Shape of data tensor: (5054, 235)
Shape of label tensor: (5054,)
Found 400000 word vectors.

Model: "sequential"
_____
Layer (type)                 Output Shape              Param #
=================================================================
embedding (Embedding)        (None, None, 100)         1000000
_____
lstm (LSTM)                  (None, 100)               80400
_____
dense (Dense)                (None, 1)                 101
=================================================================
Total params: 1,080,501
Trainable params: 1,080,501
Non-trainable params: 0
_____
None

Train on 200 samples, validate on 1000 samples
Epoch 1/10
200/200 [==============================] - 2s 10ms/sample - loss: 0.6307 - acc: 0.6350 -
    f1_m: 0.3386 - precision_m: 0.9722 - recall_m: 0.2519 - val_loss: 0.6393 - val_acc:
    0.6090 - val_f1_m: 0.6948 - val_precision_m: 0.5351 - val_recall_m: 0.9980
Epoch 2/10
```

```
200/200 [==============================] — 2s 9ms/sample — loss: 0.5955 — acc: 0.6250 —
    f1_m: 0.4697 — precision_m: 0.7757 — recall_m: 0.5645 — val_loss: 0.4988 — val_acc:
    0.8390 — val_f1_m: 0.8320 — val_precision_m: 0.7773 — val_recall_m: 0.8969
Epoch 3/10
200/200 [==============================] — 1s 7ms/sample — loss: 0.4542 — acc: 0.8700 —
    f1_m: 0.8178 — precision_m: 0.8778 — recall_m: 0.7895 — val_loss: 0.4353 — val_acc:
    0.8150 — val_f1_m: 0.8213 — val_precision_m: 0.7220 — val_recall_m: 0.9551
Epoch 4/10
200/200 [==============================] — 1s 7ms/sample — loss: 0.4383 — acc: 0.8000 —
    f1_m: 0.7469 — precision_m: 0.8571 — recall_m: 0.7411 — val_loss: 0.4922 — val_acc:
    0.7600 — val_f1_m: 0.7883 — val_precision_m: 0.6582 — val_recall_m: 0.9867
Epoch 5/10
200/200 [==============================] — 1s 7ms/sample — loss: 0.3999 — acc: 0.8250 —
    f1_m: 0.8295 — precision_m: 0.8415 — recall_m: 0.8676 — val_loss: 0.3702 — val_acc:
    0.8500 — val_f1_m: 0.8494 — val_precision_m: 0.7830 — val_recall_m: 0.9289
Epoch 6/10
200/200 [==============================] — 1s 6ms/sample — loss: 0.2866 — acc: 0.9050 —
    f1_m: 0.8928 — precision_m: 0.8759 — recall_m: 0.9155 — val_loss: 0.3480 — val_acc:
    0.8490 — val_f1_m: 0.8474 — val_precision_m: 0.7793 — val_recall_m: 0.9292
Epoch 7/10
200/200 [==============================] — 1s 7ms/sample — loss: 0.2383 — acc: 0.9050 —
    f1_m: 0.9028 — precision_m: 0.8486 — recall_m: 0.9646 — val_loss: 0.4870 — val_acc:
    0.7960 — val_f1_m: 0.7277 — val_precision_m: 0.9137 — val_recall_m: 0.6066
Epoch 8/10
200/200 [==============================] — 1s 6ms/sample — loss: 0.5132 — acc: 0.7800 —
    f1_m: 0.7525 — precision_m: 0.7843 — recall_m: 0.8197 — val_loss: 0.4105 — val_acc:
    0.8100 — val_f1_m: 0.7571 — val_precision_m: 0.8881 — val_recall_m: 0.6615
Epoch 9/10
200/200 [==============================] — 1s 6ms/sample — loss: 0.2962 — acc: 0.8850 —
    f1_m: 0.8746 — precision_m: 0.9322 — recall_m: 0.8448 — val_loss: 0.3297 — val_acc:
    0.8650 — val_f1_m: 0.8500 — val_precision_m: 0.8417 — val_recall_m: 0.8593
Epoch 10/10
200/200 [==============================] — 1s 6ms/sample — loss: 0.2154 — acc: 0.9500 —
    f1_m: 0.9505 — precision_m: 0.9344 — recall_m: 0.9677 — val_loss: 0.3243 — val_acc:
    0.8710 — val_f1_m: 0.8628 — val_precision_m: 0.8376 — val_recall_m: 0.8928
1264/1264 [==============================] — 1s 938us/sample — loss: 0.3474 — acc: 0.8536 —
    f1_m: 0.8261 — precision_m: 0.7815 — recall_m: 0.8890
Accuracy: 85.36%
Precision: 78.15%
Recall: 88.90%
F1 score: 82.61%
```

**Chapter 6**

# Conclusion

I have implemented three separate argument mining techniques, applicable to both structural and statistical approach. For the **structural approach**, there was created a dictionary of lexical cues that usually characterize argumentative structures based on (Knott and Dale, 1994). The algorithm created was tested in two IBM corpora containing only argumentative sentences, the accuracy seems to be around **15 to 17 %**. As regards the statistical approach, a corpora has to be curated in order to include both argumentative and non-argumentative instances. That is why Wikipedia articles were scrapped assuming that there are no arguments included, and the sentences gathered were cross-checked by the algorithm created for the structural approach. Every sentence that did not have argumentative cues, based on the algorithm, was added to the corpora as non-argumentative, while the others considered as ambiguous and were checked by two annotators so as to remove misleading indicators of dictionary. Afterwards, two algorithms were built by using as training data a part of this corpora. **Random forest classification algorithm** was approached by two different methods of training. The first technique was by using the tokenized sentences which leads to an accuracy of **56.83%**, precision **100%**, recall **0.55%**, f1 score **71.13%**, while the second was by using features like counter of words, uppercase and punctuation characters that have an accuracy of **79.42%**, precision **74.08%**, recall **78.72%**, and f1 score **76.33%**. Finally, the other algorithm of statistical approach that was built is the **LSTM-RNN** that concludes to an accuracy of **85.44%**, precision of **78.15%**, recall of **88.90%**, f1 score of **82.61%**.

As the results of both structural and statistical algorithms reveal, linguistic rules are not capable of successfully identifying arguments. On the other hand, learning algorithms are a better fit to argument mining, with the most accurate one to be the LSTM-RNN algorithm. One of the main goals as regarding future work is to apply the algorithms built in a fully annotated data-set found in GitHub's issues.

# References

Aharoni, Ehud et al. (2014a). "A Benchmark Dataset for Automatic Detection of Claims and Evidence". In: *Proceedings of the 25th International Conference on Computational Linguistics (COLING)*, pp. 1489–1500. URL: http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.672.6321.

— (2014b). "A Benchmark Dataset for Automatic Detection of Claims and Evidence". In: *Proceedings of the 25th International Conference on Computational Linguistics (COLING)*, pp. 1489–1500. URL: http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.672.6321.

AI, Explosion. *spaCy*. URL: https://spacy.io/.

Bar-Haim, Roy et al. (2017). "Stance Classification of Context-Dependent Claims". In: *"Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers"*. Valencia, Spain: Association for Computational Linguistics, pp. 251–261. URL: https://www.aclweb.org/anthology/E17-1024.

Biran, Or and Owen Rambow (2011). "Identifying justifications in written dialogs". In: *Proceedings - 5th IEEE International Conference on Semantic Computing, ICSC 2011* October 2011, pp. 162–168. DOI: 10.1109/ICSC.2011.41.

Budzynska, Katarzyna and Serena Villata (2015). "Argument Mining". In: *IEEE Intelligent Informatics Bulletin*.

Budzynska, Katarzyna et al. (2014). "A model for processing illocutionary structures and argumentation in debates". In: *Proceedings of the 9th International Conference on Language Resources and Evaluation, LREC 2014*, pp. 917–924.

Cabrio, Elena and Serena Villata (2012). "Natural language arguments: A combined approach". In: *Frontiers in Artificial Intelligence and Applications* 242, pp. 205–210. ISSN: 09226389. DOI: 10.3233/978-1-61499-098-7-205.

Chollet, Francois (2017). *Deep Learning with Python*. Manning Publications.

Eckle-Kohler, Judith, Roland Kluge, and Iryna Gurevych (2015). "On the role of discourse markers for discriminating claims and premises in argumentative discourse". In: *Conference Proceedings - EMNLP 2015: Conference on Empirical Methods in Natural Language Processing* September, pp. 2236–2242.

Feng, Vanessa Wei and Graeme Hirst (2011). "Classifying arguments by scheme". In: *ACL-HLT 2011 - Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies* 1, pp. 987–996.

Garcia Villalba, Maria Paz and Patrick Saint-Dizier (2012). "Some facets of argument mining for opinion analysis". In: *Frontiers in Artificial Intelligence and Applications* 245.1, pp. 23–34. ISSN: 09226389. DOI: 10.3233/978-1-61499-111-3-23.

*Glove*.

Habernalt, Ivan and Iryna Gurevych (2016). "Which argument is more convincing? Analyzing and predicting convincingness of Web arguments using bidirectional LSTM". In: *54th Annual Meeting of the Association for Computational Linguistics, ACL 2016 - Long Papers* 3, pp. 1589–1599.

John Lawrence, Chris Reed, Colin Allen, Simon McAlister, Andrew Ravenscroft, David Bourget (2014). "Mining Arguments From 19th Century Philosophical Texts Using Topic

Based Modelling". In: *Proceedings of the First Workshop on Argumentation Mining*, pp. 79–87.

Knott, Alistair and Robert Dale (1994). "Using linguistic phenomena to motivate a set of coherence relations". In: *Discourse processes* 18.1, pp. 35–62.

Lawrence, John and Chris Reed (2016). "Argument Mining Using Argumentation Scheme Structures". In: DOI: 10.3233/978-1-61499-686-6-379.

Lawrence, John and Chsris Reed (2015). "Combining Argument Mining Techniques". In: Association for Computational Linguistics (ACL), pp. 127–136. DOI: 10.3115/v1/w15-0516.

Levy, Ran et al. (2014). "Context dependent claim detection". In: *COLING 2014 - 25th International Conference on Computational Linguistics, Proceedings of COLING 2014: Technical Papers*, pp. 1489–1500.

Lippi, Marco and Paolo Torroni (2015). "Context-independent claim detection for argument mining". In: *IJCAI International Joint Conference on Artificial Intelligence*. Vol. 2015-January. International Joint Conferences on Artificial Intelligence, pp. 185–191. ISBN: 9781577357384.

— (2016). "Argumentation mining: State of the art and emerging trends". In: *ACM Transactions on Internet Technology* 16.2. ISSN: 15576051. DOI: 10.1145/2850417.

Mirkin, Shachar et al. (2017). "A Recorded Debating Dataset". In: *arXiv preprint arXiv:1709.06438*.

Mochales, Raquel and Marie Francine Moens (2011). "Argumentation mining. MARGOT: a web server for argumentation mining". In: *Artificial Intelligence and Law* 19.1, pp. 1–22. ISSN: 09248463. DOI: 10.1007/s10506-010-9104-x. arXiv: arXiv:1502.07526v1. URL: http://argumentationmining.disi.unibo.it/resources.html.

Moens, Marie Francine et al. (2007). "Automatic detection of arguments in legal texts". In: *Proceedings of the International Conference on Artificial Intelligence and Law*, pp. 225–230. ISBN: 1595936807. DOI: 10.1145/1276318.1276362.

Müller, Andreas C and Sarah Guido (2016). *Introduction to Machine Learning with Python*. Tech. rep. URL: www.wowebook.orgWOW!eBookwww.wowebook.org.

Palau, Raquel Mochales and Marie Francine Moens (2009). "Argumentation mining: The Detection, Classification and Structuring of Arguments in Text". In: *Belgian/Netherlands Artificial Intelligence Conference*, pp. 351–352. ISSN: 15687805.

Park, Joonsuk and Claire Cardie (2015). "Identifying Appropriate Support for Propositions in Online User Comments". In: pp. 29–38. DOI: 10.3115/v1/w14-2105.

Rinott, Ruty et al. (2015). "Show me your evidence - An automatic method for context dependent evidence detection". In: *Conference Proceedings - EMNLP 2015: Conference on Empirical Methods in Natural Language Processing* September, pp. 440–450.

Stab, Christian and Iryna Gurevych (2014). "Identifying argumentative discourse structures in persuasive essays". In: *EMNLP 2014 - 2014 Conference on Empirical Methods in Natural Language Processing, Proceedings of the Conference*, pp. 46–56.

Van Eemeren, Frans H, Peter Houtlosser, and AF Snoeck Henkemans (2007). *Argumentative indicators in discourse: A pragma-dialectical study*. Vol. 12. Springer Science & Business Media.

Walton, Douglas (2009). "Argumentation Theory: A Very Short Introduction". In: *Argumentation in Artificial Intelligence*.

— (2011). "Argument mining by applying argumentation schemes". In: *Studies in Logic* 4.April 2012, pp. 38–64. URL: http://www.studiesinlogic.net/english/UploadFiles{\_}1698/201104/20110415074459727.pdf.

Webber, B., M. Egg, and V. Kordoni (2012). "Discourse structure and language technology". In: *Natural Language Engineering* 18.4, pp. 437–490. ISSN: 13513249. DOI: 10.1017/S1351324911000337.

# Appendicies

# Appendicies 1

# Scrap Data from Wikipedia

```
1   import urllib.request
2   import re
3   from inscriptis import get_text
4
5
6   def wiki(theme):
7     url = "https://en.wikipedia.org/wiki/" + theme
8     html = urllib.request.urlopen(url).read().decode('utf-8')
9
10    text = get_text(html)
11
12    with open('../datasets/wiki_' + theme + '.txt', 'w') as out:
13      for row in text.split('\n'):
14        if len(row) >= 80 and not row[0].isdigit() and not row[1].isdigit() and not row[2] ==
      '*':
15          row = re.sub(r'\[\d+]', '', row)
16          row = row.rstrip('\n')
17          out.write(row)
18          out.write('\n')
19
20
21  if __name__ == '__main__':
22
23    # Video games topic
24    wiki('Video_game-related_health_problems')
25    wiki('Video_game_addiction_in_China')
26    wiki('Video_game_addiction')
27    wiki('Gaming_disorder')
28    wiki('2017_in_video_gaming')
29    wiki('2019_in_video_gaming')
30    wiki('History_of_video_games')
31    wiki('PC_game')
32    wiki('Video_game_culture')
33    wiki('Gaming_computer')
34    wiki('Video_game_console')
35    wiki('Video_game')
36    wiki('Video_gaming_in_the_United_States')
37    wiki('Video_game_music')
38    wiki('Video_game_industry')
39    wiki('Video_game_development')
40    wiki('Game_design')
41    wiki('Video_game_programmer')
42    wiki('Early_history_of_video_games')
43    wiki('Video_gaming_in_Japan')
44    wiki('Video_game_crash_of_1983')
45    wiki('Sixth_generation_of_video_game_consoles')
46    wiki('Video_gaming_in_China')
47    wiki('1980s_in_video_gaming')
48    wiki('Home_computer')
49    wiki('Nintendo')
50    wiki('Game_Boy')
51    wiki('Fourth_generation_of_video_game_consoles')
52    wiki('Game')
53    wiki('The_Game_Awards')
54
55    # Democracy topic
56    wiki('Democracy')
```

```
57    wiki('History_of_democracy')
58    wiki('Athenian_democracy')
59    wiki('Representative democracy')
60    wiki('Direct_democracy')
61    wiki('Types_of_democracy')
62    wiki('Democracy_Index')
63    wiki('Criticism_of_democracy')
64
65    # Multiculturalism
66    wiki('Multiculturalism')
67    wiki('Criticism_of_multiculturalism')
68    wiki('Cultural_pluralism')
69    wiki('Multiculturalism_in_Canada')
70    wiki('Multicultural_education')
71    wiki('Multicultural_and_diversity_management')
72    wiki('Multiculturalism_in_Australia')
73    wiki('Multicultural_transruption')
74
```

# Appendicies 2

# Remove Duplicate Sentences

```python
import os

FILE_PATH = os.path.abspath(os.path.dirname(__file__))


def remove_duplicate_rows(file_name):
  sentences = []

  with open(os.path.join(FILE_PATH, '../Results/' + file_name), mode='r') as txt_file:
    reader = txt_file.read()

    for sentence in reader.split("\n"):
      sentences.append(sentence)

  with open ( os.path.join ( FILE_PATH, '../Results/new_' + file_name ), mode='w' ) as
    txt_file_new:
      for i in range(len(sentences)):
        if i == 0:
          txt_file_new.write(sentences[i])
          txt_file_new.write("\n")
        elif i < len(sentences) and sentences[i] != sentences[i-1]:
          txt_file_new.write(sentences[i])
          txt_file_new.write("\n")


if __name__ == '__main__':
  remove_duplicate_rows('CDEdata.csv')
```

# Appendicies 3

# Concat Results of checked Data

```python
import os
import glob
import csv

FILE_PATH = os.path.abspath(os.path.dirname(__file__))


def combine_results():

  csvfiles = glob.glob(FILE_PATH + '/../Results/checked/*.csv')
  dataset = csv.writer(open(FILE_PATH + '/../Results/dataset.csv', 'w'), delimiter=',')

  for files in csvfiles:
    rd = csv.reader(open(files, 'r'))

    for row in rd:
      if 'wiki' in files:
        if "False" in row:
          dataset.writerow(row)
      else:
        if "True" or "False" in row:
          dataset.writerow(row)


def find_ambiguous():

  csvfiles = glob.glob(FILE_PATH + '/../Results/checked/*.csv')
  found_ambiguous = csv.writer(open(FILE_PATH + '/../Results/found_ambiguous.csv', 'w'),
    delimiter=',')

  for files in csvfiles:
    rd = csv.reader(open(files, 'r'))

    for row in rd:
      if 'wiki' in files:
        if "True" in row:
          found_ambiguous.writerow(row)


if __name__ == '__main__':
  combine_results()
  find_ambiguous()
```

# Appendicies 4

# Gather Ambiguous Sentences found in Wiki articles

## 4.1 ../Results/found_ambiguous_results.csv

In the **terminal** type:

```
$python3 getAmbiguousSentences.py > ../Results/found_ambiguous_results.csv
```

## 4.2 getAmbiguousSentences.py

```python
import json
from configParser import choose_function
from getArguments import spaCy, pos_tagged, check_regex

"""
Description: by using a dictionary that includes words often used
in arguments, it is identified if given sentences are
arguments or not. Part of speech tagging from spaCy is used
for this purpose as well and it is imported by getArguments.py file.
"""


def check_dictionary(doc, dictionary):
  """None
  This function checks if any of the words in the sentence exists
   in the dictionary given. If it does, then it is checked if
   this word's part of speech match with its value given in
   dictionary. If they match, then the word is added in a
   list named keyword_found.

  :param doc: pos tagged sentence from spacy function
  :param dictionary: dictionary that has as keywords words
    and as value their part of speech
  :return: keywords found in the given sentence
  """

  keyword_found = []

  for key, value in dictionary.items():
    if check_regex(doc, key) is not None:
      if len(value) == 1:
        for key2 in value:
          if checz_regex(doc, key + ' ' + key2) is not None
            and check_regex(doc, key2) is not None and
            pos_tagged(doc, check_regex(doc, key2).text) is not 'None':
            if value[key2] == pos_tagged(doc, check_regex(doc, key2).text)[1] or \
              value[key2] == pos_tagged(doc, check_regex(doc, key2).text)[2]:
              keyword_found.append(key + " + " + key2)
      elif value == 'none':
        keyword_found.append(key)
      elif len(value) != 1 and check_regex(doc, value)
        is not None:
```

```python
43            keyword_found.append(key)
44        elif pos_tagged(doc, key)[1] == value or \
45          pos_tagged(doc, key)[2] == value:
46            keyword_found.append(key)
47
48    return keyword_found
49
50
51 if __name__ == '__main__':
52
53    with open('../dict/dictionary.json', 'r') as dict:
54      dictionary = json.load(dict)
55
56    sentences = choose_function("found_ambiguous.csv")
57
58    if sentences != 'No dataset found':
59      for sentence in sentences:
60
61        print('"' + str(sentence[0]).strip('b') + '",' +
62          str(check_dictionary(spaCy(sentence[0]), dictionary)))
63
64    else:
65      print(sentences)
```

# Appendicies 5

# Enumerate ambiguous sentences

```python
import os
import csv
import json

FILE_PATH = os.path.abspath(os.path.dirname(__file__))


def enumerate_ambiguous(file_name):
  with open('../dict/dictionary.json', 'r') as dict:
    dictionary = json.load(dict)

    with open(os.path.join(FILE_PATH, '../Results/' + file_name), mode='r', encoding='utf-8
    ') as file:
      reader = csv.reader(file)
      counters = {}

      count = 0
      count_rows = 0

      for row in reader:
        firstArgumentWordFlag = False

        count_rows += 1
        # first cell includes sentences, not keywords
        for column in row[1:]:

          if firstArgumentWordFlag == False:
            testStr = "['"
            if column.startswith(testStr):
            firstArgumentWordFlag = True
            else:
              continue

          column = column.strip("['] ")
          column = column.strip('"')

          if column in dictionary:
            count += 1
            if column in counters:
              counters[column] += 1
            else:
              counters[column] = 1
          else:
            if ' + ' in column:
              count += 1
              firstPart = column.split(" + ")[0]

              if firstPart in counters:
                counters[firstPart] += 1
              else:
                counters[firstPart] = 0
            else:
              print(row)
              print('Not found: ', column)

    save_array_to_csv(counters)
```

```
57
58  def save_array_to_csv(my_dict):
59    with open(FILE_PATH + '/../Results/found_Ambiguous_keywords.csv', 'w', newline='') as f:
60      w = csv.writer(f)
61      w.writerows(my_dict.items())
62
63
64  if __name__ == '__main__':
65    enumerate_ambiguous('found_ambiguous_results.csv')
```

# Appendicies 6

# Extract keywords from pdf

```python
"""
Description: extract words which are often used in arguments
(based on a paper),and create a dictionary based on these words
(key of the dict) and their specific, if they have one, part
of speech (value of the dict) in arguments
"""

import json
import textract
import os
import csv

FILE_PATH = os.path.abspath(os.path.dirname(__file__))  # path of this file


def extract_data():
  """
  By using textract library, this function extracts the whole pdf file

  pdf: paper called 'Using Linguistic Phenomena to Motivate a Set
  of Coherence Relations'
  """

  text = textract.process(os.path.join(FILE_PATH,
  "../Reading/cues-UsingLinguisticPhenomenaMotivateCoherenceRelations_Knott93.pdf"))
  save(text)


def save(text):
  """
  This function saves extracted text to a csv file
  """
  if not os.path.exists("../dict" ):
    os.makedirs("../dict")

  with open(os.path.join(FILE_PATH, "../dict/data.csv"),
   mode='wb') as csv_file:
    csv_file.write(text)

  modify_csv_file("../dict/data.csv")  # modify extracted text


def modify_csv_file(data):
  """
  This function modifies csv file in order to keep those words
  we are interested in
  """

  flag = 0

  with open(os.path.join(FILE_PATH, data)) as inp:
    reader = csv.reader(inp)

    with open(os.path.join(FILE_PATH, "../dict/data2.csv"),
      mode='w') as out:
        for row in reader:
          if len(row) > 0 and row[0] == "Phrase":
```

```
58              flag = 1
59              continue
60          if len(row) == 0 or row[0].isdigit():
61              flag = 0
62          if flag == 1 and len(row) > 0:
63              out.write(row[0])
64              out.write("\n")
65  check_words("../dict/data2.csv", "../dict/data.csv")


def check_words(data2, data):
    """
    This function adds or removes words that considered as useful
    or not
    """

    exclude_words = ['after', 'and', 'as soon as', 'before',
            'at first', 'at first sight', 'earlier',
            'fisrt of all', 'for', 'inasmuch as',
            'later', 'much sooner', 'not because',
            'now', 'if not', 'if so',
            'in the beginning', 'in the end',
            'in the meantime', 'in turn',
            'much later', 'not', 'notwithstanding that',
            'suppose', 'the more often', 'this time',
            'presumably because', 'when', 'where',
            'previously', 'regardless of that', 'rather',
            'after that', 'as', 'simply because', 'then',
            'true', 'until', 'again', 'and/or', 'or',
            'else', 'even']

    include_words = ['for the reason that', 'besides',
            '(E|e)(ither).+?(or)','(N|n)(either).+?(nor)',
            'in one hand', 'in this case',  'on one side',
            'as a matter of fact', 'in point of fact',
            'presumably', 'provided that',
            'regardless', 'rather than', 'simply',
            'as an example', 'in addition']

    test_words = {'even though': 'none', 'first': 'adv',
            'against': 'none', 'last': 'adv',
            'more': {'[a-z]*ly': 'adv'},
            'most': {'[a-z]*ly': 'adv'}, 'if': 'none',
            '(T|t)(he more).+?(the more)': 'none',
            '(T|t)(he more).+?(the less)': 'none',
            'naturally': 'none', 'once again': 'none',
            'once more': 'none', 'surely': 'none',
            'second': 'adv', 'so': 'mark', 'third': 'adv',
            'too': '(too)($|[\.])', 'should say': 'none',
            'might say': 'none', 'may say': 'none',
            'could say': 'none', 'while': 'mark',
            'as a start': 'none', 'in order to': 'none',
            'in order that': 'none', 'still': 'adv',
            'that is': 'none', 'since': 'mark',
            'yet': '(Y|y)(et)[^\.].', 'that': 'mark'}

    with open(os.path.join(FILE_PATH, data2), 'r') as inp, \
        open(os.path.join(FILE_PATH, data), 'w') as out:

        for row in csv.reader(inp):
            if row[0] in exclude_words:
                continue
            else:
                out.write(row[0])
                out.write("\n")

        for word in include_words:
            out.write(word)
            out.write("\n")

    create_dictionary("../dict/data.csv", test_words)
```

```
130
131  def create_dictionary(data, test_words):
132      """
133      This function creates a .json file that includes a dictionary of
134      the words from the csv file created before and some additional
135      words for testing
136      """
137
138      dictionary = test_words
139
140      with open(os.path.join(FILE_PATH, data), 'r') as inp:
141
142        for row in csv.reader(inp):
143          if "\x05" in row[0]:
144            row[0] = row[0].replace('\x05', 'fi')   # correct words from pdf extraction
145
146          if row[0] in test_words.keys():
147            continue
148          else:
149            dictionary.update({row[0]: 'none'})
150
151      with open('../dict/dictionary.json', 'w') as dict:
152        json.dump(dictionary, dict)
153
154
155  if __name__ == '__main__':
156      extract_data()
```

# Appendicies 7

# Load Data-sets based on their configurations

```python
1  from xlrd import open_workbook
2  import configparser
3  import os
4  import re
5  import csv
6
7  FILE_PATH = os.path.abspath(os.path.dirname(__file__))
8
9
10 def py23_str(value):
11   """
12   This function tries to convert a string to unicode. Because
13   of the fact that this conversion differ from python 3
14   to python 2, here are checked both possibilities so as
15   the program to run in both python 3 and 2.
16
17   :param value: sentence to be converted from string to unicode
18   :return: converted input
19   """
20
21   try:  # Python 2
22     return unicode(value, errors='ignore', encoding='utf-8')
23   except NameError:  # Python 3
24     try:
25       return str(value, errors='ignore', encoding='utf-8')
26     except TypeError:  # Wasn't a bytes object, no need to decode
27       return str(value)
28
29
30 def get_sentences_csv(dataset_number):
31   """
32   This function reads files with .csv extension
33
34   :dataset_number: number that refers to order (starts from 0)
35   of a dataset in datasets.ini
36
37   :return: a list of sentences
38   """
39   sentences = []
40
41   path, _, column, is_argument = get_parameters_dataset(dataset_number)
42
43   with open(os.path.join(FILE_PATH, path), mode='r') as dataset:
44     reader = csv.reader(dataset)
45     for sentence in reader:
46       if is_argument is not None:
47         sentences.append([str(sentence[int(column)]), str(sentence[int(is_argument)])])
48       else:
49         sentences.append([str(sentence[int(column)]), 'True'])
50
51   sentences.pop(0)
52   return sentences
53
```

```
54
55  def get_sentences_xls(dataset_number):
56      """
57      This function reads files with .xls extension
58
59      :dataset_number: number that refers to order (starts from 0)
60      of a dataset in datasets.ini
61      :return: a list of sentences
62      """
63      sentences = []
64
65      path, sheet, column, is_argument = get_parameters_dataset(dataset_number)
66
67      reader = open_workbook(path, on_demand=True)
68      sheet = reader.sheet_by_name(sheet)
69      if is_argument is not None:
70          for cell, cell2 in zip(sheet.col(int(column)), sheet.col(int(is_argument))):
71              sentences.append([cell.value.encode("utf-8"), cell2.value.encode("utf-8")])
72      else:
73          for cell in sheet.col(int(column)):
74              sentences.append([cell.value.encode("utf-8"), 'True'])
75
76      sentences.pop(0)
77      return sentences
78
79
80  def get_sentences_txt(dataset_number):
81      """
82      This function reads files with .txt or none extension
83
84      :dataset_number: number that refers to order (starts from 0)
85              of a dataset in datasets.ini
86
87      :return: a list of sentences
88      """
89      sentences = []
90
91      path, _,  _, _ = get_parameters_dataset(dataset_number)
92
93      with open(os.path.join(FILE_PATH, path), mode='r') as txt_file:
94          reader = txt_file.read()
95
96          for sentence in reader.split('.'):
97              sentences.append([sentence])
98
99      return sentences
100
101
102 def get_parameters_dataset(dataset):
103      """
104      This function gets the arguments of a specific dataset from datasets.ini
105
106      :dataset: number that refers to order (starts from 0)
107              of a dataset or the name of dataset in datasets.ini
108      :return: section['path'] + file_name: path of dataset
109      sheet: sheet that data are in it if it is an .xls file
110      column: column of sentences to be identified as arguments or not
111      is_argument: column which reveals if a specific sentence is
112              an argument or not
113      """
114      dataset_number, config = check_validity_of_dataset(dataset)
115
116      section = config.sections()[dataset_number]  # each section is a name of a file with data
117      section = config[section]
118      file_name = re.match(r".*: (.*)>", str(section), re.MULTILINE)
119      file_name = file_name.group(1)
120
121      try:
122          sheet = section['sheet']
123      except KeyError:
124          sheet = None
125
```

```
126    try:
127        is_argument = section['is_argument']
128    except KeyError:
129        is_argument = None
130
131    try:
132        column = section['column']
133    except KeyError:
134        column = None
135
136    return section['path'] + file_name, sheet, column, is_argument
137
138
139 def check_validity_of_dataset(dataset):
140    """
141    This function checks of a dataset exists in dataset.ini or not
142
143    :dataset: number that refers to order (starts from 0) of
144    a dataset or the name of dataset in datasets.ini
145    :return: dataset_number: returns the order of given
146    dataset in datasets.ini config: returns object config
147    from datasets.ini
148    """
149    config = configparser.ConfigParser()
150    config.read('../datasets/datasets.ini')
151
152    if dataset in config:
153        dataset_number = config.sections().index(dataset)
154    elif dataset < len(config.sections()):
155        dataset_number = dataset
156
157    return dataset_number, config
158
159
160 def choose_function(dataset):
161    """
162    This function checks the extension of a datasets and chooses
163    an appropriate method to read the file
164
165    :dataset: number that refers to order (starts from 0)
166    of a dataset or the name of dataset in datasets.ini
167    :return: a list of sentences if dataset exits
168    otherwise 'No dataset found'
169    """
170
171    try:
172        dataset_number = int(check_validity_of_dataset(dataset)[0])
173
174        try:
175            _, extension = dataset.rsplit('.', 1)
176        except ValueError:
177            extension = None
178
179        if extension == 'xls':
180            return get_sentences_xls(dataset_number)
181        elif extension == 'csv':
182            return get_sentences_csv(dataset_number)
183        elif extension == 'txt' or extension is None:
184            return get_sentences_txt(dataset_number)
185
186    except TypeError:
187        return 'No dataset found'
```

# Appendicies 8

# Find Argumentative Sentences

```
1  import json
2  import spacy
3  from __future__ import division
4  from configParser import choose_function, os, py23_str, re
5
6  """
7  Description: by using a dictionary that includes words often used
8  in arguments, this file identifies if given sentences are
9  arguments or not. Part of speech tagging from spaCy is used
10 for this purpose as well.
11 """
12
13 FILE_PATH = os.path.abspath(os.path.dirname(__file__))
14
15 tp = 0
16 tn = 0
17 fp = 0
18 fn = 0
19
20
21 def spaCy(sentence):
22   """
23   By using spaCy, this function gets a sentence and returns every
24   word's part of speech
25
26   :param sentence: input to be tokenized
27   :return: tokenized sentence
28   """
29
30   nlp = spacy.load('en')
31   doc = nlp(py23_str(sentence))
32
33   return doc
34
35
36 def pos_tagged(doc, word):
37   """
38   This function gets a tagged sentence from spaCy and a specific
39   word and return its part of speech and its dependency
40
41   :param doc: pos tagged sentence from spacy function
42   :param word: a word that we are interested to learn its part of
43     speech
44   :return: word, its part of speech(pos) and its dependence in the
45     given sentence or None
46   """
47
48   word = word.lower()
49
50   for token in doc:
51     if token.text.lower() == word:
52       return [token.text, token.pos_.lower(), token.dep_.lower()]
53
54   return 'None'
55
56
57 def check_regex(doc_regex, regex):
```

```
58     """
59     By using spaCy's function called match, this function is
60     checking if a specific regular expression is represented
61     by a given sentence
62
63     :param doc_regex: pos tagged sentence from spacy function
64     :param regex: a regular expression
65     :return: the part of the sentence that is indicated in the given
66       regex otherwise None
67     """
68
69     regex = re.compile(r''+regex)
70
71     for match in re.finditer(regex, doc_regex.text.lower()):
72       start, end = match.span()  # get matched indices
73       word_found = doc_regex.char_span(start, end)  # create Span from indices
74
75       return word_found
76
77     return None
78
79
80  def check_dictionary(doc, dictionary):
81     """
82     This function checks if any of the words in the sentence exists
83      in the dictionary given. If it does, then it is checked if
84      this word's part of speech match with its value given in
85      dictionary. If they match, then the word is added in a
86      list named keyword_found.
87
88     :param doc: pos tagged sentence from spacy function
89     :param dictionary: dictionary that has as keywords words
90       and as value their part of speech
91     :return: True if the list keyword_found is not empty or False
92          if it is empty
93     """
94
95     keyword_found = []
96
97     for key, value in dictionary.items():
98       if check_regex(doc, key) is not None:
99         if len(value) == 1:
100          for key2 in value:
101            if checz_regex(doc, key + ' ' + key2) is not None
102              and check_regex(doc, key2) is not None and
103              pos_tagged(doc, check_regex(doc, key2).text) is not 'None':
104               if value[key2] == pos_tagged(doc, check_regex(doc, key2).text)[1] or \
105                 value[key2] == pos_tagged(doc, check_regex(doc, key2).text)[2]:
106                 keyword_found.append(key + " + " + key2)
107        elif value == 'none':
108          keyword_found.append(key)
109        elif len(value) != 1 and check_regex(doc, value)
110         is not None:
111          keyword_found.append(key)
112        elif pos_tagged(doc, key)[1] == value or \
113         pos_tagged(doc, key)[2] == value:
114          keyword_found.append(key)
115
116    if len(keyword_found) != 0:
117      return 'True'
118    else:
119      return 'False'
120
121
122 def check_validity(real_value, given_value):
123    """
124    This function checks if two given values match
125
126    :param real_value:  real value is the result of check_dictionary
127      function
128    :param given_value: given value is the value given by analysts
129      into dataset
```

```
130      :return: Correct results if they match or Wrong results if they
131        do not match
132      """
133      global tn, tp, fn, fp
134
135      if real_value in given_value:
136        if real_value == 'True':
137          tp = tp + 1
138        else:
139          tn = tn + 1
140        return "Correct results"
141      else:
142        if real_value == 'False':
143          fp = fp + 1
144        else:
145          fn = fn + 1
146        return "Wrong results"
147
148
149  def precision():
150      if (tp + fp) != 0:
151        return tp/(tp + fp)
152      else:
153        return "Integer division by zero"
154
155
156  def recall():
157      if (tp + fn) != 0:
158        return tp/(tp + fn)
159      else:
160        return "Integer division by zero"
161
162
163  def f1_score(precision, recall):
164      if (precision + recall) != 0:
165        return 2 * (precision * recall) / (precision + recall)
166      else:
167        return "Integer division by zero"
168
169
170  if __name__ == '__main__':
171
172      with open('../dict/dictionary.json', 'r') as dict:
173        dictionary = json.load(dict)
174
175      sentences = choose_function("found_fp.csv")
176      labeled_data = False
177
178      if sentences != 'No dataset found':
179            for sentence in sentences:
180
181          if len(sentence) == 2:
182            print('"' + str(sentence[0]).strip('b') + '",' + 'True') # for csv
183            # print ( str ( sentence[0] ).strip ( 'b' ) + ',' + 'True' )
184            labeled_data = True
185
186          else:
187            print('"' + str(sentence[0]).strip('b') + '",' +
188              check_dictionary(spaCy(sentence[0]), dictionary)) # for csv
189
190      else:
191        print(sentences)
```

# Appendicies 9

# Statistical Approach: Random Forest algorithm

```python
1   import pandas as pd
2   import seaborn as sn
3   from sklearn.preprocessing import OneHotEncoder
4   from sklearn.model_selection import train_test_split
5   from sklearn.ensemble import RandomForestClassifier
6   from sklearn.metrics import confusion_matrix
7   import matplotlib.pyplot as plt
8   import string
9   import csv
10  import os
11  import numpy as np
12
13
14  """
15  Description: implementing Random Forest classifier model
16  to an annotated corpora that contains
17  both argumentative and none sentences
18  """
19
20  FILE_PATH = os.path.abspath(os.path.dirname(__file__))  # path of this file
21
22
23  def load_dataset(path):
24      """
25      Loading dataset of a given path, and creating features based on the sentences given.
26      The first feature is a counter of words included in each sentence,
27      the second feature is a counter of uppercase characters, while
28      the third feature is a counter of special characters (punction)
29      :param path: full path of dataset
30      :return: a list of sentences and their labels, and a set of features
31      """
32
33      # load & prepare data
34      with open(path, 'r') as file:
35          dataset = csv.reader(file, delimiter=",")
36          df = pd.DataFrame(dataset)
37
38      del df[2]  # column full of non labels = df[1]
39      df[1] = [len(sentences.split()) for sentences in df[0]]  # Word Count'
40      df[2] = [sum(char.isupper() for char in sentence) \
41          for sentence in df[0]]  # 'Uppercase Char Count'
42
43      df[3] = [sum(char in string.punctuation for char in sentence) \
44          for sentence in df[0]]  # 'Special Char Count'
45      df[4] = pd.factorize(labels)[0]  # switch False to 0 and True to 1
46
47      return df
48
49
50  def tokenize_sentences(df):
51      """
52      Tokenizing raw sentences by using One-hot encoding
53      into a format that a computer can understand
```

```python
54
55      :param df: the dataframe that includes 4 columns
56             [sentences, word Counter, uppercase counter,
57             special char counter, label]
58      :return: tokenized sentences, labels
59      """
60      vectorizer = OneHotEncoder(handle_unknown='ignore')
61      vectorizer.fit([[line.strip()] for line in df[0]])
62      sentences = vectorizer.transform([[line.strip()] for line in df[0]]).toarray()
63
64      # sentences = vectorizer.fit_transform(df[0])
65      labels = df[4]
66      return sentences, labels
67
68
69  def define_model(x_train, y_train):
70      """
71      This function defines and trains the Random Forest model
72
73      :param x_train: the training sentences
74      :param y_train: the sentences' labels
75      :return: trained model
76      """
77      model = RandomForestClassifier()
78      model.fit(x_train, y_train)
79      return model
80
81
82  def precision(tp, fp):
83      if (tp + fp) != 0:
84          return tp/(tp + fp)
85      else:
86          return "Integer division by zero"
87
88
89  def recall(tp, fn):
90      if (tp + fn) != 0:
91          return tp/(tp + fn)
92      else:
93          return "Integer division by zero"
94
95
96  def f1_score(precision, recall):
97      if type(recall) != str and type(precision) != str and (precision + recall) != 0:
98          return 2 * (precision * recall) / (precision + recall)
99      else:
100         return "Integer division by zero"
101
102
103 def testing_model(model, x_test, y_test):
104     """
105     This function evaluates the model on the test set
106
107     :param model: trained model
108     :param x_test: the testing set of sentences
109     :param y_test: the sentences' labels
110     :return: a matrix of the predicted and real values
111     """
112     score = model.score(x_test, y_test)
113     print("Accuracy: %.2f%%" % (score * 100))
114
115     y_predicted = model.predict(x_test)
116     cm = confusion_matrix(y_test, y_predicted)
117
118     fp = cm[0][1]
119     fn = cm[1][0]
120     tp = cm[1][1]
121
122     prec = precision(tp, fp)
123     rec = recall(tp, fn)
124
125     print("Precision: %.2f%%" % (prec * 100))
```

```
126    print("Recall: %.2f%%" % (rec * 100))
127    print("F1 Score: %.2f%%" % (f1_score(prec, rec) * 100) )
128
129    print(cm)
130    return cm
131
132
133 def plot_results(cm):
134    """
135    This function is showing a heatmap plot based on
136    the values of the confusion matrix that contains
137    the real an predicted values.
138    :param cm: matrix of both real and predicted values
139    """
140    plt.figure(figsize=(10, 7))
141    sn.heatmap(cm, annot=True)
142    plt.xlabel('Prediction')
143    plt.ylabel('Truth')
144    plt.show()
145
146
147 if __name__ == '__main__':
148    df = load_dataset('../Results/dataset.csv')
149    sentences, labels = tokenize_sentences(df)
150
151    #  using features instead of the sentences
152    #
153    features = np.asarray(df[df.columns[1:4]].values)
154    # split dataset into test and train data
155    x_train, x_test, y_train, y_test = \
156      train_test_split(features, labels, test_size=0.33)
157
158    model = define_model(x_train, y_train)
159    cm = testing_model(model, x_test, y_test)
160    plot_results(cm)
161
162    #  using sentences without their features
163    #
164    # split dataset into test and train data
165    x_train, x_test, y_train, y_test = train_test_split(sentences, labels, test_size=0.33)
166
167    model = define_model(x_train, y_train)
168    cm = testing_model(model, x_test, y_test)
169    plot_results(cm)
```

# Appendicies 10

# Statistical Approach: LSTM-RNN algorithm

```python
1   import csv
2   import os
3   import pandas as pd
4   import numpy as np
5   import matplotlib.pyplot as plt
6   from sklearn.model_selection import train_test_split
7   from tensorflow.python.keras.preprocessing.text import Tokenizer
8   from tensorflow.python.keras.preprocessing.sequence \
9     import pad_sequences
10  from tensorflow.python.keras.layers import Dense, LSTM, Embedding
11  from tensorflow.python.keras.models import Sequential
12  from keras import backend
13
14  """
15  Description: implement lstm model to an annotated corpora
16  that contains both argumentative and none sentences
17  """
18
19  # path of this file
20  FILE_PATH = os.path.abspath(os.path.dirname(__file__))
21
22
23  def load_dataset(path):
24    """
25    Loading dataset of a given path
26    :param path: full path of dataset
27    :return: a list of sentences and their labels
28    """
29
30    # load & prepare data
31    with open(path, 'r') as file:
32      dataset = csv.reader(file, delimiter=",")
33      df = pd.DataFrame(dataset)
34
35    del df[2]  # column full of none
36
37    df[1] = pd.factorize(df[1])[0]   # False switched to 0 and True to 1
38    texts = list(df[0].values)
39    labels = list(df[1])
40
41    return texts, labels
42
43
44  def get_tokenized_text(max_words, texts_train):
45    """
46    Tokenizing raw sentences.
47    :param max_words: considering only the top max_words of dataset
48            (most frequent ones)
49    :param texts_train: sentences to for tokenization
50    :return: tokenizer, sequences, word_index
51    """
52    tokenizer = Tokenizer(num_words=max_words)
53    tokenizer.fit_on_texts(texts_train)
```

```
54      sequences = tokenizer.texts_to_sequences(texts_train)
55      word_index = tokenizer.word_index
56
57      print('Found %s unique tokens.' % len(word_index))
58      return tokenizer, sequences, word_index
59
60
61  def pad_sentences(sequences, labels_train):
62      """
63      This function transforms a list of num_sentences sequences
64      (lists of integers) into a 2D Numpy array of shape
65      (num_sentnece, num_timesteps).
66      The num_timesteps is either the maxlen argument if provided,
67      or the length of the longest sequence otherwise.
68      Sequences that are shorter than num_timesteps are padded
69      with value at the end.
70
71      :param sequences: list of tokenized sentences
72      :param labels_train: list of sequences' labels
73      :return: a 2D Numpy array of sentences and labels
74      """
75      data = pad_sequences(sequences)
76      labels_train = np.asarray(labels_train)
77      print('Shape of data tensor:', data.shape)
78      print('Shape of label tensor:', labels_train.shape)
79
80      # Splits the data into a training set and a validation set,
81      # but first shuffles the data
82      indices = np.arange(data.shape[0])
83      np.random.shuffle(indices)
84      data = data[indices]
85      labels_train = labels_train[indices]
86
87      return data, labels_train
88
89
90  def word_embedding(glove_file, max_words, word_index):
91      """
92      Parsing the GloVe's word-embeddings file in order
93      to build an index that maps words (as strings)
94      to their vector representation (as number vectors),
95      and then build an embedding matrix that will
96      be loaded into an Embedding layer later on.
97
98      :param glove_file: path of Glove's word-embedding
99      :param max_words: considering only the top max_words
100              of dataset (most frequent ones)
101      :param word_index: index of words after tokenization
102      :return: embedding_dim, embedding_matrix
103          of shape (max_words, embedding_dim)
104      """
105
106      embeddings_index = {}
107      f = open(os.path.join(FILE_PATH, glove_file))
108      for line in f:
109          values = line.split()
110          word = values[0]
111          coefs = np.asarray(values[1:], dtype='float32')
112          embeddings_index[word] = coefs
113      f.close()
114      print('Found %s word vectors.' % len(embeddings_index))
115
116      embedding_dim = 100
117      embedding_matrix = np.zeros((max_words, embedding_dim))
118      for word, i in word_index.items():
119          if i < max_words:
120              embedding_vector = embeddings_index.get(word)
121              if embedding_vector is not None:
122                  # Words not found in the embedding
123                  # index will be all zeros.
124                  embedding_matrix[i] = embedding_vector
125
```

```
126    return embedding_dim, embedding_matrix
127
128
129  def lstm_model(max_words, embedding_dim, embedding_matrix):
130    """
131    This function defines the model and
132    loads pretrained word embedding into the Embedding layer
133
134    :param max_words: considering only the top max_words
135              of dataset (most frequent ones)
136    :param embedding_dim: number of embedding dimensions used
137    :param embedding_matrix: a matrix of shape shape
138                (max_words, embedding_dim)
139    :return: the defined model
140    """
141
142    # define model
143    model = Sequential()
144    model.add(Embedding(max_words, embedding_dim))
145    model.add(LSTM(100))
146    model.add(Dense(1, activation='sigmoid'))
147    print(model.summary())
148
149    # Loading pretrained word embeddings into the Embedding layer
150    model.layers[0].set_weights([embedding_matrix])
151    model.layers[0].trainable = False
152
153    return model
154
155
156  def recall_m(y_true, y_pred):
157    """
158    Calculating recall metric
159
160    :param y_true: the real value of a sentence
161    :param y_pred: the estimated value of a sentence
162    :return: recall metric
163    """
164    true_positives = backend.sum(backend.round(backend.clip(y_true * y_pred, 0, 1)))
165    possible_positives = backend.sum(backend.round(backend.clip(y_true, 0, 1)))
166    recall = true_positives / (possible_positives + backend.epsilon())
167    return recall
168
169
170  def precision_m(y_true, y_pred):
171    """
172    Calculating precision metric
173
174    :param y_true: the real value of a sentence
175    :param y_pred: the estimated value of a sentence
176    :return: precision metric
177    """
178    true_positives = backend.sum(backend.round(backend.clip(y_true * y_pred, 0, 1)))
179    predicted_positives = backend.sum(backend.round(backend.clip(y_pred, 0, 1)))
180    precision = true_positives / (predicted_positives + backend.epsilon())
181    return precision
182
183
184  def f1_m(y_true, y_pred):
185    """
186    Calculating F1 Score metric
187
188    :param y_true: the real value of a sentence
189    :param y_pred: the estimated value of a sentence
190    :return: f1_score metric
191    """
192    precision = precision_m(y_true, y_pred)
193    recall = recall_m(y_true, y_pred)
194    return 2 * ((precision * recall) / (precision + recall + backend.epsilon()))
195
196
197  def training_model(x_train, y_train, x_val, y_val, model):
```

```python
198    """
199    This function is responsible for compiling and
200    training the model
201
202    :param x_train: sentences for training model (200 sentences)
203    :param y_train: labels of training sentences
204    :param x_val: sentences for validating trained model
205            (1000 sentences)
206    :param y_val: labels of validation sentences
207    :param model: defined model
208    :return: historing of training, trained model
209    """
210
211    model.compile(loss='binary_crossentropy', optimizer='rmsprop',
212            metrics=['accuracy'])
213    history = model.fit(x_train, y_train, epochs=10, batch_size=128,
214            validation_data=(x_val, y_val))
215    model.save_weights('pre_trained_glove_model.h5')
216
217    return history, model
218
219
220 def plotting_results(history):
221    """
222    This function is showing plots of the training and
223    validation results based on accuracy and loss metrics.
224
225    :param history: history of training procedure
226    """
227    acc = history.history['acc']
228    val_acc = history.history['val_acc']
229    loss = history.history['loss']
230    val_loss = history.history['val_loss']
231    epochs = range(1, len(acc) + 1)
232
233    plt.plot(epochs, acc, 'bo', label='Training acc')
234    plt.plot(epochs, val_acc, 'b', label='Validation acc')
235    plt.title('Training and validation accuracy')
236    plt.legend()
237    plt.figure()
238
239    plt.plot(epochs, loss, 'bo', label='Training loss')
240    plt.plot(epochs, val_loss, 'b', label='Validation loss')
241    plt.title('Training and validation loss')
242    plt.legend()
243    plt.show()
244
245
246 def testing_model(texts_test, labels_test, tokenizer, model):
247    """
248    This function tokenizes the data of test set,
249    and evaluates the model on the test set
250
251    :param texts_test: sentences for testing model
252    :param labels_test: labels of the testing sentences
253    :param tokenizer: tokenizer
254    :param model: trained model
255    :return: accuracy of the created model
256    """
257
258    # Tokenizing the data of the test set
259    sequences = tokenizer.texts_to_sequences(texts_test)
260    x_test = pad_sequences(sequences)
261    y_test = np.asarray(labels_test)
262
263    # Evaluating the model on the test set
264    model.load_weights('pre_trained_glove_model.h5')
265    loss, accuracy, f1_score, precision, recall = model.evaluate(x_test, y_test)
266
267    print("Accuracy: %.2f%%" % (accuracy * 100))
268    print("Precision: %.2f%%" % (precision * 100))
269    print("Recall: %.2f%%" % (recall * 100))
```

```
270    print("F1 score: %.2f%%" % (f1_score * 100))
271
272
273  def main():
274    texts, labels = load_dataset('../Results/dataset.csv')
275    texts_train, texts_test, labels_train, labels_test = \
276      train_test_split(texts, labels, test_size=0.2)
277
278    max_words = 10000  # Considers only the top 10,000 words
279                # in the dataset
280    training_samples = 200
281    validation_samples = 1000\usepackage{xcolor}
282    \definecolor{maroon}{cmyk}{0, 0.87, 0.68, 0.32}
283    \definecolor{halfgray}{gray}{0.55}
284    \definecolor{ipython_frame}{RGB}{207, 207, 207}
285    \definecolor{ipython_bg}{RGB}{247, 247, 247}
286    \definecolor{ipython_red}{RGB}{186, 33, 33}
287    \definecolor{ipython_green}{RGB}{0, 128, 0}
288    \definecolor{ipython_cyan}{RGB}{64, 128, 128}
289    \definecolor{ipython_purple}{RGB}{170, 34, 255}
290
291    tokenizer, sequences, word_index = \
292      get_tokenized_text(max_words, texts_train)
293    data, labels_train = pad_sentences(sequences, labels_train)
294
295    x_train = data[:training_samples]
296    y_train = labels_train[:training_samples]
297    x_val = \
298      data[training_samples: training_samples + validation_samples]
299    y_val = \
300      labels_train[training_samples: training_samples + validation_samples]
301
302    embedding_dim, embedding_matrix \
303    = word_embedding('../Reading/glove.6B/glove.6B.100d.txt', \
304            max_words, word_index)
305
306    model = lstm_model(max_words, embedding_dim, embedding_matrix)
307    history, model = training_model(x_train, y_train, \
308                  x_val, y_val, model)
309    plotting_results(history)
310    testing_model(texts_test, labels_test, tokenizer, model)
311
312
313  if __name__ == '__main__':
314    main()
```