

ATHENS UNIVERSITY OF ECONOMICS AND BUSINESS

BACHELOR RESEARCH ASSIGNMENT

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 26, 2019

Contents

1	Introduction	2
1.1	Definition	2
1.2	Research Goal	2
1.3	Assignment's Structure	3
2	State-of-the-Art	4
3	Methods	7
3.1	Structural Approach	7
3.2	Machine Learning Approach	8
3.2.1	randomForest.py	8
3.2.2	LSTM	8
4	Data Curation	10
4.1	Data Used	10
5	Results	14
5.1	Structural Approach	14
5.2	Machine Learning Approach	15
6	Conclusion	16
6.1	Future Work	16
	Appendix	17
1	Scrap Data from Wikipedia	19
2	Remove Duplicate Sentences	21
3	Concat Results of checked Data	22
4	False-Positives Creation	24
4.1	../Results/found_fp_results.csv	24
4.2	getFalsePositives.py	24
5	Enumerate false positives	26
6	Extract keywords from pdf	28
7	Extract Keywords from pdf	32
8	Find Argumentative Sentences	37
9	Machine Learning Approach: Random Forest	42

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-term 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 in hand coded rules and the machine learning approach.

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 machine learning 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.

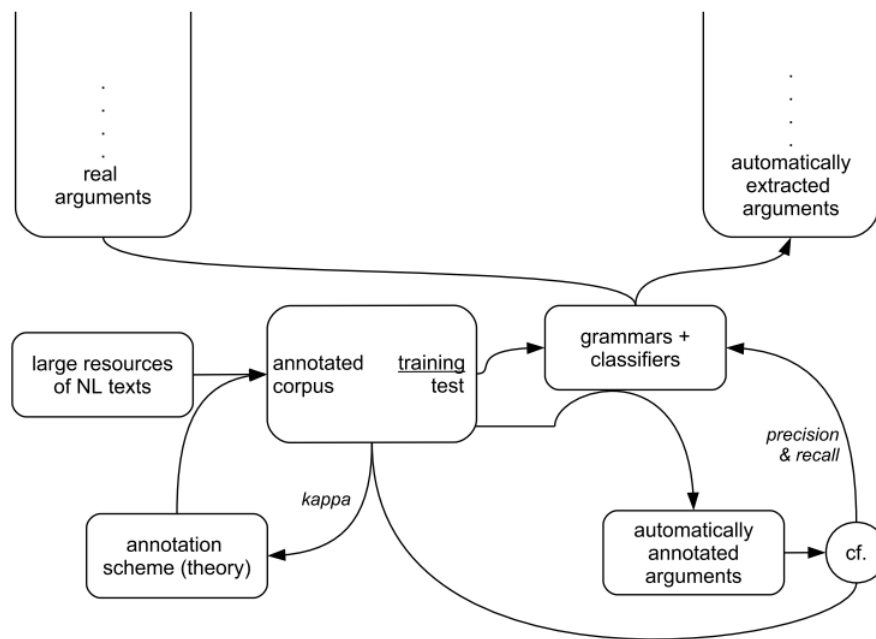


FIGURE 2.1: Natural language processing techniques

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 (Habernal 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; Ecker-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; Ecker-Kohler, Kluge, and Gurevych, 2015), Maximum Entropy classifiers (Mochales and Moens, 2011), and Decision Trees and Random Forests (Stab and Gurevych, 2014; Ecker-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
Ecker-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)[^\\.\].
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 encountering precision, recall and *f1_score*, that are metrics for reviewing algorithm's results. These metrics will be presented in more details at Chapter 5.

3.2 Machine Learning Approach

3.2.1 *randomForest.py*

the Random Forest algorithm introduced a robust, practical take on decision-tree learning that involves building a large number of specialized decision trees and then ensembling their outputs. Random forests are applicable to a wide range of problems you could say that they're almost always the second-best algorithm for any shallow machine-learning task. Chollet, 2017

3.2.2 LSTM

It can be understood as either a sequence of characters or a sequence of words, but its most common to work at the level of words. The deep-learning sequence-processing models introduced in the following sections can use text to produce a basic form of natural-language understanding, sufficient for applications including document classification, sentiment analysis, author identification, and even question-answering (QA) (in a constrained context). Of course, keep in mind throughout this chapter that none of these deep-learning models truly understand text in a human sense; rather, these models can map the statistical structure of written language, which is sufficient to solve many simple textual tasks. Deep learning for natural-language processing is pattern recognition applied to words, sentences, and paragraphs, in much the same way that computer vision is pattern recognition applied to pixels. Chollet, 2017

Collectively, the different units into which you can break down text (words, characters, or n-grams) are called tokens, and breaking text into such tokens is called tokenization. All text-vectorization processes consist of applying some tokenization scheme and then associating numeric vectors with the generated tokens. These vectors, packed into sequence tensors, are fed into deep neural networks. There are multiple ways to associate a vector with a token. In this section, I'll present two major ones: one-hot encoding of tokens, and token embedding (typically used exclusively for words, and called word embedding). Chollet, 2017

Note that Keras has built-in utilities for doing one-hot encoding of text at the word level or character level, starting from raw text data. You should use these utilities, because they take care of a number of important features such as stripping special characters from strings and only taking into account the N most common words in your dataset (a common restriction, to avoid dealing with very large input vector spaces). Chollet, 2017

A variant of one-hot encoding is the so-called one-hot hashing trick, which you can use when the number of unique tokens in your vocabulary is too large to handle explicitly. Instead of explicitly assigning an index to each word and keeping a reference of these indices in a dictionary, you can hash words into vectors of fixed size. This is typically done with a very lightweight hashing function. The main advantage of this method is that it does away with maintaining an explicit word index, which saves memory and allows online encoding of the data (you can generate token vectors right away, before you've seen all of the available data). The one drawback of this approach is that it's susceptible to hash collisions: two different words may end up with the same hash, and subsequently any machine-learning model looking at these hashes won't be able to tell the difference between these words. The likelihood of hash collisions decreases when the dimensionality of the hashing space is much larger than the total number of unique tokens being hashed. Chollet, 2017

Another popular and powerful way to associate a vector with a word is the use of dense word vectors, also called word embeddings. Whereas the vectors obtained through one-hot encoding are binary, sparse (mostly made of zeros), and very high-dimensional (same dimensionality as the number of words in the vocabulary), word embeddings are low-dimensional floating-point vectors (that is, dense vectors, as opposed to sparse vectors); see figure 6.2. Unlike the word vectors obtained via one-hot encoding, word embeddings are learned from data. It's common to see word embeddings that are 256-dimensional, 512-dimensional, or 1,024-dimensional when dealing with very large vocabularies. On the other hand, one-hot encoding words generally leads to vectors that are 20,000-dimensional or greater (capturing a vocabulary of 20,000 tokens, in this case). So, word embeddings pack more information into far fewer dimensions. Chollet, 2017

the geometric relationships between word vectors should reflect the semantic relationships between these words. Word embeddings are meant to map human language into a geometric space. For instance, in a reasonable embedding space, you would expect synonyms to be embedded into similar word vectors; and in general, you would expect the geometric distance (such as L2 distance) between any two word vectors to relate to the semantic distance between the associated words (words meaning different things are embedded at points far away from each other, whereas related words are closer). In addition to distance, you may want specific directions in the embedding space to be meaningful. To make this clearer, let's look at a concrete example. Chollet, 2017

It's thus reasonable to learn a new embedding space with every new task. Fortunately, backpropagation makes this easy, and Keras makes it even easier. It's about learning the weights of a layer: the Embedding layer. Chollet, 2017

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.

4.1 Data Used

Most of the argumentative sentences included in the our corpora were found on two IBM data-sets created for this purpose (Table 4.1), including 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., 2014	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 already existing annotated data-sets referred to 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, which makes it impossible to train a supervised algorithm classifier in identifying 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 IBM curated data. 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 this sentence 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 are more specifically about Video games, Democracy and Multiculturalism. The code used for the scraping process is 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., 2014	wiki12_articles/Gender_representation_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 that is aligned with Appendix 6. The sentences classified as non-arguments were added to the created data-set, while the others were considered as controversial (Table 4.3).

Topic of Wikipedia	Total Data	Not Used	Controversial sentences	Non-Arguments
Early History 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), and the corpora that will be used in machine learning algorithms is finally composed of 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: Argumentative and No-Argumentative Data Used

Additionally, a file that includes all the controversial sentences was created (Appendix 4). This file indicates all the keywords responsible for these controversial results. The indicators- described in previous Chapter and found in Wikipedia articles- were counted using code of Appendix 5, and the following table was the outcome of that enumeration.

either	53	for example	93
(E e)(ither).+?(or)	44	as a result	991
also	456	as well	1065
because	121	notably	13
for instance	17	but	326
that is	582	while	197
actually	26	against	51
still	85	so that	558
though	87	besides	4
furthermore	66	eventually	41
clearly	5	if	107
since	34	on the grounds that	527
although	170	third	6
consequently	4	as a consequence	946
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	488
in order to	43	just as	877
at the same time	7	of course	2
finally	13	as an example	846
first	29	even though	95
lastly	823	equally	7
whereas	828	(T t)(he more).+?(the more)	2
regardless	7	for this reason	803
simply because	109	by contrast	781
naturally	4	at any rate	1
in short	2	in this case	727
such that	398	essentially	5
(N n)(either).+?(nor)	41	whenever	4
second	5	as long as	609
even then	3	as a matter of fact	578
accordingly	3	provided that	283
conversely	3	alternatively	3
afterwards	6	thereafter	1
meanwhile	127	once again	3
once more	1	above all	1
by comparison	7	surely	2
undoubtedly	2	on the one side	1
at first	8	presumably	2
after all	1	what is more	1
certainly	1	anyway	1

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 sentences in Wikipedia. A representative number of them was examined by the authors of this research papers, and the keywords that did not pointed out argumentative statements and removed from the dictionary

are represented in the following table. (+ table)

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 a group of annotated data were executed.

A set of metrics, recall, precision and F1 score, were used in order to evaluate the indicators selected for recognizing argumentative sentences. These metrics are using four counters, 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, while false negative (fn) how many times human identified a sentence as argumentative while algorithm did not.

- **Precision** indicates a metric of correctly identified instances:

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

- **Recall** measures the times algorithm missed out an instance:

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

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

$$P = \frac{2 * (R * P)}{R + P}$$

The results based on these three metrics are presented in the table bellow, and they were applied to ten corpora in total. It has to be mentioned that the data-sets were manually annotated by this paper's authors, and included a small amount of both argumentative and non-argumentative sentences.

Source Data	file	Precision	Recall	F1 Score
Mirkin et al., 2017	asr/DJ_1_ban-video-games_pro.wav.asr.txt	1.0	1.0	1.0
Mirkin et al., 2017	asr/EH_1_ban-video-games_pro.wav.asr.txt	1.0	1.0	1.0
Mirkin et al., 2017	asr/HE_1_ban-video-games_pro.wav.asr.txt	0.8	0.727	0.762
Mirkin et al., 2017	asr/SN_1_video-games_pro.wav.asr.txt	1.0	0.928	0.963
Mirkin et al., 2017	asr/TL_1_ban-video-games_pro.wav.asr.txt	0.727	0.888	0.799
Mirkin et al., 2017	asr/YB_1_ban-video-games_pro.wav.asr.txt	0.555	0.833	0.667
Aharoni et al., 2014	wiki12_articles/Gender_representation_in_video_games	0.167	0.4	0.235

5.2 Machine Learning Approach

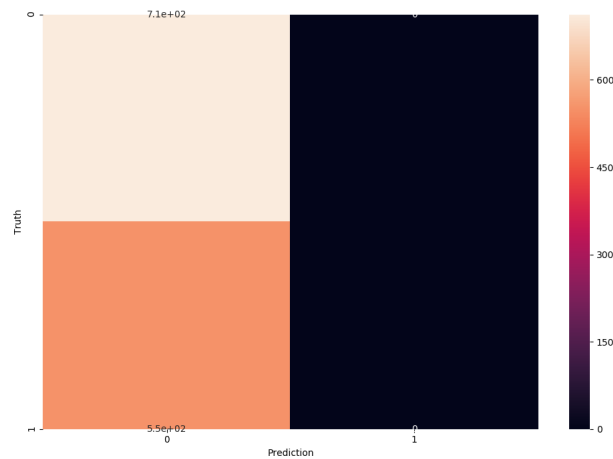


FIGURE 5.1: Results of Random Forest Algorithm in the created dataset

Chapter 6

Conclusion

6.1 Future Work

Appendix

Appendix 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].
15                isdigit() and not row[2] == '*':
16                row = re.sub(r'\[\d+\]', '', row)
17                row = row.rstrip('\n')
18                out.write(row)
19                out.write('\n')
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')
```

Appendix 2

Remove Duplicate Sentences

```

1 import os
2
3 FILE_PATH = os.path.abspath(os.path.dirname(__file__))
4
5
6 def remove_duplicate_rows(file_name):
7     sentences = []
8
9     with open(os.path.join(FILE_PATH, '../Results/' + file_name),
10               mode='r') as txt_file:
11         reader = txt_file.read()
12
13         for sentence in reader.split("\n"):
14             sentences.append(sentence)
15
16     with open ( os.path.join ( FILE_PATH, '../Results/new_' +
17                               file_name ), mode='w' ) as txt_file_new:
18         for i in range(len(sentences)):
19             if i == 0:
20                 txt_file_new.write(sentences[i])
21                 txt_file_new.write("\n")
22             elif i < len(sentences) and sentences[i] != sentences[i-1]:
23                 txt_file_new.write(sentences[i])
24                 txt_file_new.write("\n")
25
26 if __name__ == '__main__':
27     remove_duplicate_rows('CDEdata.csv')

```


Appendix 3

Concat Results of checked Data

```
1 import os
2 import glob
3 import csv
4
5 FILE_PATH = os.path.abspath(os.path.dirname(__file__))
6
7
8 def combine_results():
9
10     csvfiles = glob.glob(FILE_PATH + '/../Results/checked/*.csv')
11     dataset = csv.writer(open(FILE_PATH + '/../Results/dataset.csv',
12                               'w'), delimiter=',')
13
14     for files in csvfiles:
15         rd = csv.reader(open(files, 'r'))
16
17         for row in rd:
18             if 'wiki' in files:
19                 if "False" in row:
20                     dataset.writerow(row)
21             else:
22                 if "True" or "False" in row:
23                     dataset.writerow(row)
24
25 def find_false_positives():
26
27     csvfiles = glob.glob(FILE_PATH + '/../Results/checked/*.csv')
28     found_fps = csv.writer(open(FILE_PATH + '/../Results/found_fp.
29                                csv', 'w'), delimiter=',')
30
31     for files in csvfiles:
32         rd = csv.reader(open(files, 'r'))
33
34         for row in rd:
35             if 'wiki' in files:
36                 if "True" in row:
37                     found_fps.writerow(row)
38
```

```
39 if __name__ == '__main__':  
40     combine_results()  
41     find_false_positives()
```

Appendix 4

False-Positives Creation

4.1 ../Results/found_fp_results.csv

In the **terminal** type:

```
1 python3 getFalsePositives.py > ../Results/found_fp_results.csv
```

4.2 getFalsePositives.py

```
1 import json
2 from configparser import choose_function
3 from getArguments import spaCy, pos_tagged, check_regex
4
5 """
6 Description: by using a dictionary that includes words often used
7 in arguments, it is identified if given sentences are
8 arguments or not. Part of speech tagging from spaCy is used
9 for this purpose as well and it is imported by getArguments.py
10 file.
11 """
12
13 def check_dictionary(doc, dictionary):
14     """None
15     This function checks if any of the words in the sentence exists
16     in the dictionary given. If it does, then it is checked if
17     this word's part of speech match with its value given in
18     dictionary. If they match, then the word is added in a
19     list named keyword_found.
20
21     :param doc: pos tagged sentence from spacy function
22     :param dictionary: dictionary that has as keywords words
23         and as value their part of speech
24     :return: keywords found in the given sentence
25     """
26
27     keyword_found = []
28
29     for key, value in dictionary.items():
30         if check_regex(doc, key) is not None:
```

```
31     if len(value) == 1:
32         for key2 in value:
33             if checzech_regex(doc, key + '_' + key2) is not None
34             and check_regex(doc, key2) is not None and
35             pos_tagged(doc, check_regex(doc, key2).text) is not '
                None':
36                 if value[key2] == pos_tagged(doc, check_regex(doc,
                    key2).text)[1] or \
37                 value[key2] == pos_tagged(doc, check_regex(doc, key2)
                    .text)[2]:
38                     keyword_found.append(key + "_+_" + key2)
39     elif value == 'none':
40         keyword_found.append(key)
41     elif len(value) != 1 and check_regex(doc, value)
42     is not None:
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_fp.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)
```

Appendix 5

Enumerate false positives

```

1 import os
2 import csv
3 import json
4
5 FILE_PATH = os.path.abspath(os.path.dirname(__file__))
6
7
8 def enumerate_false_positives(file_name):
9     with open('../dict/dictionary.json', 'r') as dict:
10         dictionary = json.load(dict)
11
12     with open(os.path.join(FILE_PATH, '../Results/' + file_name),
13               mode='r') as file:
14         reader = csv.reader(file)
15         counters = []
16
17         for row in reader:
18             # first cell includes sentences, not keywords
19             for column in row[1:]:
20                 print(column)
21                 column = column.strip("[']_")
22                 column = column.strip('"')
23
24                 if column in dictionary:
25                     flag = True
26
27                     for counter in counters:
28                         if column in counter[0]:
29                             counter[1] += 1
30                             flag = False
31
32                     if flag:
33                         counters.append([column, 1])
34
35     save_array_to_csv(counters)
36
37 def save_array_to_csv(arrayOfArrays):
38     file = csv.writer(open(FILE_PATH +

```

```
39     '../Results/found_FP_keywords.csv', 'w'), delimiter=',')
40
41     for array in arrayOfArrays:
42         file.writerow([array[0]] + [array[1]])
43
44
45 if __name__ == '__main__':
46     enumerate_false_positives('found_fp_results.csv')
```

Appendix 6

Extract keywords from pdf

```

1  """
2  Description: extract words which are often used in arguments
3  (based on a paper), and create a dictionary based on these words
4  (key of the dict) and their specific, if they have one, part
5  of speech (value of the dict) in arguments
6  """
7
8  import json
9  import textract
10 import os
11 import csv
12
13 FILE_PATH = os.path.abspath(os.path.dirname(__file__)) # path of
    this file
14
15
16 def extract_data():
17     """
18     By using textract library, this function extracts the whole pdf
19     file
20     pdf: paper called 'Using Linguistic Phenomena to Motivate a Set
21     of Coherence Relations'
22     """
23
24     text = textract.process(os.path.join(FILE_PATH,
25     "../Reading/cues-
        UsingLinguisticPhenomenaMotivateCoherenceRelations_Knott93.
        pdf"))
26     save(text)
27
28
29 def save(text):
30     """
31     This function saves extracted text to a csv file
32     """
33     if not os.path.exists("../dict"):
34         os.makedirs("../dict")
35

```

```

36 with open(os.path.join(FILE_PATH, "../dict/data.csv"),
37           mode='wb') as csv_file:
38     csv_file.write(text)
39
40     modify_csv_file("../dict/data.csv") # modify extracted text
41
42
43 def modify_csv_file(data):
44     """
45     This function modifies csv file in order to keep those words
46     we are interested in
47     """
48
49     flag = 0
50
51     with open(os.path.join(FILE_PATH, data)) as inp:
52         reader = csv.reader(inp)
53
54         with open(os.path.join(FILE_PATH, "../dict/data2.csv"),
55                   mode='w') as out:
56             for row in reader:
57                 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")
66
67
68 def check_words(data2, data):
69     """
70     This function adds or removes words that considered as useful
71     or not
72     """
73
74     exclude_words = ['after', 'and', 'as_soon_as', 'before',
75                      'at_first', 'at_first_sight', 'earlier',
76                      'fisrt_of_all', 'for', 'inasmuch_as',
77                      'later', 'much_sooner', 'not_because',
78                      'now', 'if_not', 'if_so',
79                      'in_the_beginning', 'in_the_end',
80                      'in_the_meantime', 'in_turn',
81                      'much_later', 'not', 'notwithstanding_that',
82                      'suppose', 'the_more_often', 'this_time',
83                      'presumably_because', 'when', 'where',
84                      'previously', 'regardless_of_that', 'rather',
85                      'after_that', 'as', 'simply_because', 'then',
86                      'true', 'until', 'again', 'and/or', 'or',

```



```

87         'else', 'even']
88
89     include_words = ['for_the_reason_that', 'besides',
90                     '(E|e)(ither).+?(or)', '(N|n)(either).+?(nor)',
91                     'in_one_hand', 'in_this_case', 'on_one_side',
92                     'as_a_matter_of_fact', 'in_point_of_fact',
93                     'presumably', 'provided_that',
94                     'regardless', 'rather_than', 'simply',
95                     'as_an_example', 'in_addition']
96
97     test_words = {'even_though': 'none', 'first': 'adv',
98                  'against': 'none', 'last': 'adv',
99                  'more': {'[a-z]*ly': 'adv'},
100                  'most': {'[a-z]*ly': 'adv'}, 'if': 'none',
101                  '(T|t)(he_more).+?(the_more)': 'none',
102                  '(T|t)(he_more).+?(the_less)': 'none',
103                  'naturally': 'none', 'once_again': 'none',
104                  'once_more': 'none', 'surely': 'none',
105                  'second': 'adv', 'so': 'mark', 'third': 'adv',
106                  'too': '(too)($|[\.\.])', 'should_say': 'none',
107                  'might_say': 'none', 'may_say': 'none',
108                  'could_say': 'none', 'while': 'mark',
109                  'as_a_start': 'none', 'in_order_to': 'none',
110                  'in_order_that': 'none', 'still': 'adv',
111                  'that_is': 'none', 'since': 'mark',
112                  'yet': '(Y|y)(et)[^\.\.]', 'that': 'mark'}
113
114     with open(os.path.join(FILE_PATH, data2), 'r') as inp, \
115           open(os.path.join(FILE_PATH, data), 'w') as out:
116
117         for row in csv.reader(inp):
118             if row[0] in exclude_words:
119                 continue
120             else:
121                 out.write(row[0])
122                 out.write("\n")
123
124         for word in include_words:
125             out.write(word)
126             out.write("\n")
127
128     create_dictionary("../dict/data.csv", test_words)
129
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
145             # from pdf extraction
146
147         if row[0] in test_words.keys():
148             continue
149         else:
150             dictionary.update({row[0]: 'none'})
151
152 with open('../dict/dictionary.json', 'w') as dict:
153     json.dump(dictionary, dict)
154
155 if __name__ == '__main__':
156     extract_data()
```

Appendix 7

Extract Keywords from pdf

```

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
105     datasets.ini
106
107     :dataset: number that refers to order (starts from 0)
108     of a dataset or the name of dataset in datasets.ini
109     :return: section['path'] + file_name: path of dataset
110     sheet: sheet that data are in it if it is an .xls file
111     column: column of sentences to be identified as arguments or not
112     is_argument: column which reveals if a specific sentence is
113     an argument or not
114     """
115
116     dataset_number, config = check_validity_of_dataset(dataset)
117
118     section = config.sections()[dataset_number] # each section is a
119     name of a file with data
120
121     section = config[section]
122     file_name = re.match(r".*:_(.*)>", str(section), re.MULTILINE)
123     file_name = file_name.group(1)
124
125     try:
126         sheet = section['sheet']
127     except KeyError:
128         sheet = None
129
130     try:
131         is_argument = section['is_argument']
132     except KeyError:
133         is_argument = None
134
135     try:
136         column = section['column']
137     except KeyError:
138         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 if 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 exists
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'
```

Appendix 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
74         from indices
75
76         return word_found
77
78     return None
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 check_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 '
104                         None':
105                         if value[key2] == pos_tagged(doc, check_regex(doc,
106                             key2).text)[1] or \
107                             value[key2] == pos_tagged(doc, check_regex(doc, key2)
108                                 .text)[2]:
109                             keyword_found.append(key + "_+" + key2)
110                     elif value == 'none':
111                         keyword_found.append(key)
112                     elif len(value) != 1 and check_regex(doc, value)
113                     is not None:
114                         keyword_found.append(key)
115                     elif pos_tagged(doc, key)[1] == value or \
116                     pos_tagged(doc, key)[2] == value:
117                         keyword_found.append(key)
118
119     if len(keyword_found) != 0:
120         return 'True'
121     else:
122         return 'False'
123
124 def check_validity(real_value, given_value):
125     """
126     This function checks if two given values match
127
128     :param real_value: real value is the result of check_dictionary
129                       function
130     :param given_value: given value is the value given by analysts
131                       into dataset
132     :return: Correct results if they match or Wrong results if they
133             do not match
134     """
135     global tn, tp, fn, fp
136
137     if real_value in given_value:
138         if real_value == 'True':
139             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') #
183                     # print ( str ( sentence[0] ).strip ( 'b' ) + ', ' + 'True'
184                         )
185                 labeled_data = True
186
187     else :

```

```
187         print('"' + str(sentence[0]).strip('b') + '", ' +  
188             check_dictionary(spacy(sentence[0]), dictionary)) # for  
189                 csv  
190     else:  
191         print(sentences)
```

Appendix 9

Machine Learning Approach: Random Forest

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.metrics import confusion_matrix
4 from sklearn.preprocessing import OneHotEncoder
5 import seaborn as sn
6 import pandas as pd
7 import matplotlib.pyplot as plt
8 import csv
9
10
11 def encode_dataset(sentences):
12     enc = OneHotEncoder(handle_unknown='ignore')
13     print(sentences)
14     enc.fit([[line.strip()] for line in sentences])
15     enc.categories_
16     print(enc.categories_)
17     result = enc.transform([[line.strip()] for line in sentences]).
18         toarray()
19     return result
20
21 def my_dataset():
22
23     with open('../Results/dataset.csv', 'r') as dataset:
24         dataset = csv.reader(dataset, delimiter=",")
25         # for row in dataset:
26         #     print(row)
27         df = pd.DataFrame(dataset)
28         del df[2]
29         y = pd.factorize(df[1])[0]
30         df[1] = y
31         target = df[1]
32         df[0] = encode_dataset(df[0])
33
34         x_train, x_test, y_train, y_test = train_test_split(df.drop
35             ([1], axis='columns'),

```

```
36     print(x_train , y_train)
37     print(x_test , y_test)
38     model = RandomForestClassifier(n_estimators=20)
39     model.fit(x_train , y_train)
40     print(model.score(x_test , y_test))
41
42     y_predicted = model.predict(x_test)
43     cm = confusion_matrix(y_test , y_predicted)
44     print(cm)
45
46     plt.figure(figsize=(10, 7))
47     sn.heatmap(cm, annot=True)
48     plt.xlabel('Prediction')
49     plt.ylabel('Truth')
50     plt.show()
51
52
53 if __name__ == '__main__':
54     my_dataset()
55     # False 0 -True 1
```

References

- Aharoni, Ehud et al. (2014). "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](https://doi.org/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](https://doi.org/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](https://doi.org/10.3233/978-1-61499-111-3-23).
- 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 Chrsis Reed (2015). "Combining Argument Mining Techniques". In: Association for Computational Linguistics (ACL), pp. 127–136. DOI: [10.3115/v1/w15-0516](https://doi.org/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](https://doi.org/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](https://doi.org/10.1007/s10506-010-9104-x). arXiv: [arXiv:1502.07526v1](https://arxiv.org/abs/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](https://doi.org/10.1145/1276318.1276362).
- 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](https://doi.org/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](https://doi.org/10.1017/S1351324911000337).