



Athens University of Economics & Business
Department of Management Science & Technology
MSc in Business Analytics
Machine Learning & Content Analytics
“Argument Detection Project”

<i>MSc Student_1:</i>	<i>MSc Student_2:</i>	<i>MSc Student_3:</i>
<i>Agapiou Marios</i>	<i>Kontos Christos</i>	<i>Skarlis Vasileios</i>
<i>A.M.: f2821901</i>	<i>A.M.: f2821906</i>	<i>A.M.: f2821912</i>

Table of Contents

1. Project Description	1
1.1 Objectives and business workflows of the project	2
1.2 Structure of the project	3
2. Mission	4
2.1 Previous efforts	5
2.2 Approach of the problem	6
3. Data	7
3.1 Web scrapping	7
3.2 Annotation procedure using the CLARIN·EL	7
3.3 Data Parsing	9
3.4 Pre-processing procedures	12
4. Methodology	14
4.1 Models architectures and hyperparameters	14
1. Feed Forward Neural Network	14
2. Convolutional Neural Network	18
3. Long Short-Term Memory Neural Network	21
4.2 Model comparison	24
5. Results	25
5.1 Main findings	25
5.2 Blind test set predictions	29
6. Conclusions and future work	31
Appendix	32
i) Members/Roles	32
ii) Time Plan	33
Bibliography	34

Abbreviation List

BoW	Bag of Words
CLARIN	Common Language Resources and Technology Infrastructure
CNN	Convolutional Neural Network
DESA	Department of Economic and Social Affairs
FFNN	Feed Forward Neural Network
HMM	Hidden Markov Model
MLP	Multilayer Perceptron
NLP	Natural Language Processing
LSTM	Long Short-Term Memory neural network
RNN	Recurrent Neural Network
SDG	Sustainable Development Goals
SVM	Support Vector Machines
UN	United Nations

List of Figures

Figure 1: Structure of a Feed Forward Neural Network.....	14
Figure 2: Architecture of feed-forward neural network	16
Figure 3: Structure of a Convolutional Neural Network for NLP	18
Figure 4: Architecture of Convolutional neural network	19
Figure 5: Structure of a Long Short-Term Memory Neural Network	21
Figure 6: Architecture of Bidirectional LSTM model	22
Figure 7: Confusion Matrix plot and selected metrics	27

List of Plots

Plot 1: Distribution of the data.....	10
Plot 2: Distribution of the most common characteristics of the sentences.....	11
Plot 3: 25 most common words	11
Plot 4: Division of training, validation and test datasets.....	12
Plot 5: Distribution of the maximal length of the sentences	13
Plot 6: Learning curves of the MLP model.....	25
Plot 7: Learning curves of the CNN model	26
Plot 8: Learning curves of the RNN model	26
Plot 9: Test accuracies of all models using the cross-validation procedure.....	27
Plot 10: Confusion matrix of the MLP model	28
Plot 11: Confusion matrix of the CNN model	28
Plot 12: Confusion matrix of the LSTM model	29

List of Tables

Table 1: Structure of the data	9
Table 2: Main statistical results of the data	9
Table 3: Hyperparameters of the Feed Forward Neural Network	15
Table 4: Hyperparameters of the Convolutional Neural Network.....	18
Table 5: Hyperparameters of the Long Short-Term Memory Neural Network.....	22
Table 6: Classification report of the MLP model	28
Table 7: Classification report of the CNN model	28
Table 8: Classification report of the LSTM model	29
Table 9: Blind test set predictions according to the best model.....	30

Ensure healthy lives and promote well-being for all at all ages – Argument detection using deep Neural Networks

Abstract

The aim of this project is to accurately classify and predict text content and specifically categorize sentences with claims and evidences by its relevance to the measurable indicators of the United Nations' Sustainable Development Goal #3 (Good Health and Well-Being) using a variety of Deep Neural Network models for multiclass text classification, that have been demonstrated through the years to be capable of achieving remarkable performance in sentence and document modelling. This approach we believe that will help governments to take better decisions and make more accurate actions concerning the elimination of the specific problem, which has plagued for many years millions of people. We will consider three network models: the first one is a Feed Forward Neural Network (MLP), the second is a Convolutional Neural Network (CNN) and the third one is a Long Short-Term Memory Recurrent Neural Network (LSTM). We will take into account several open access records and scientific articles collections / documents from PubMed and we will investigate on the behavior of the neural networks by defining which model give us the best results and subsequently evaluate the predictions of the best model using the PubMed scientific articles collection as test case.

1. Project Description

Ensuring healthy lives and promoting well-being for all at all ages is important to building prosperous societies. Health for all people, all over the world, is an important part of sustainable development. However, despite great strides in improving people's health and well-being in recent years, inequalities in health care access still persist. More than six million children still die before their fifth birthday each year, and only half of all women in developing regions have access to the health care they need. Various diseases exist that cause serious health issues, including epidemics like tuberculosis, HIV/AIDS, polio and malaria, where fear and discrimination limit people's ability to receive the services they need to live healthy and productive lives. Access to good health and well-being is a human right, and that is why the Sustainable Development Agenda offers a new chance to ensure that everyone can access the highest standards of health and health care— not just the wealthiest.

The Sustainable Development Goals (SDGs) or Global Goals are a collection of 17 interlinked goals designed to be a “blueprint to achieve a better and sustainable future for all”. The SDGs, set in 2015 by the United Nations General Assembly and intended to be achieved by the year 2030, are part of a UN Resolution called the “2030 Agenda”. Between the many resolutions, speeches, reports and other documents that are produced each year, the United Nations is awash in text. In particular, there is an interest in measuring how the work of the United Nations system aligns with the Sustainable Development Goals (SDGs). In particular, there is a need for a scalable, objective and consistent way to measure and classify the main arguments that are being detected in publications concerning the specific SDG target, with the ultimate goal being the insurance of healthy lives and the promotion of well-being for all at all ages.

With the increasing availability of electronic documents and the extraction of open data from any source and form, the task of automatic categorization of documents became the key method for organizing the information and knowledge discovery. Their proper classification and categorization, need text mining, machine learning and natural language processing techniques to get meaningful knowledge. In general, text classification plays an important role in information extraction and summarization, text retrieval, and question-answering. Nowadays, deep learning-based models have surpassed classical machine learning based approaches in various text classification tasks, including sentiment analysis, news categorization, question answering and natural language inference.

Therefore, the purpose of this project is to classify text content by its relevance and more specifically to accurately predict sentences including claims and evidences of various documents to the measurable indicators of the United Nations' Sustainable Development Goal #3 – Health and Well-Being (SDG3) using several Deep Neural Networks for label and multi-class text classification in order to help us identify patterns and gaining important insights, by eliminating as much as possible the time and the volume of documents one's has to read in order to acquire a specific opinion and take a corresponding action. Specifically, we developed three network models for our calculations. A Feed Forward Neural Network (MLP model), a Convolutional Neural

Network (CNN model) and a Long Short-Term Memory Recurrent Neural Network (LSTM or RNN model), using data derived from PubMed open access records and scientific articles collections, filtered for studies, using targets and indicators of the specific SDG. Finally, for the implementation of our analysis and the deployment of our models we split our initial dataset into train and validation sets using a stratified shuffle split and a cross-validation approach, by setting also each time, the appropriate hyperparameters, layers and functions at each specific model.

1.1 Objectives and business workflows of the project

The United Nations is a source of big data in the form of text. Between the many resolutions, speeches, meetings, conferences, studies, reports and internal regulations that exist and that are produced each year, the UN is awash in text. Even in a single department of the United Nations Secretariat, the amount of publications is significant. In the Department of Economic and Social Affairs (DESA), publications are central to its overall mission to support international cooperation in the pursuit of sustainable development for all. They inform development policies, global standards and norms on a wide range of development issues that affect peoples' lives and livelihoods. However, very few people are in a position to see much more than a small sliver of specialized text. Even fewer can parse the various streams into a coherent and useful picture. What is needed is a quick and objective way to analyze large quantities of United Nations publications according to a desired, criteria namely the Sustainable Development Goals (SDGs). Each SDG has several "targets", or social outcomes that the UN hopes to achieve by 2030. Each of these targets is measured using a set of indicators. These indicators represent the quantitative measurements that will be used to judge whether each SDG target has been achieved or not by 2030. Specifically, SDG3 has 14 targets and 27 indicators.

Our challenge was to gather as much as possible information from all these publications and articles referring to the specific SDG's targets and indicators, in order to predict and classify properly, claims and evidences that arise from all these different publications contained inside the PubMed open access records and scientific articles collections library. By doing this and using appropriate neural network models, we succeeded to reduce the extra noise appeared in all of these documents and successfully isolated only the important information needed, in order to get more specific insights later with the right categorization of claims and evidences drawn from the specific SDG's targets and indicators. Using machine learning models and especially deep neural network algorithms to analyze digital texts has many advantages. Over the past several years, the global Data Science community has watched the rise and steady penetration of such concepts as neural networks, Deep Learning, and back propagation. Algorithms can be used at scale with objectivity and can help identify patterns across publications and over time. This approach can also serve as a tool to explore and discover new texts and to inform the direction of future research. More importantly, this method hopefully will inspire other efforts to use modern data analytics to better understand the body of work of the United Nations. Finally, we believe that this method will help institutions as well as governments too, taking better decisions and implementing more accurate actions, concerning this SDG, in order to confront this problem and hopefully to the elimination of this problem that has plagued for many years millions of people.

1.2 Structure of the project

This project is organized as follows. Following this description, the project discusses the objectives and the necessary business workflows in which our project is integrated. In the second section is extensively analyzed previous works concerning the specific approach as well as our way of approaching the problem. The third section explains all the procedures of how we gain our data, the pre-process procedures that we follow in order to bring them into a remarkable form to develop our algorithms properly, as well as important statistic results arising from the data used in this project. The fourth section presents the methodologies and architectures that we followed for the development of our neural network models, as well as all the parameters, hyper-parameters, evaluation measures and optimizers that we used. A fifth section presents also all the results that emerged from our analyses, as well as the selection and evaluation of the best model, by including also all the necessary information gained from the predictions and the important metrics (confusion matrices, classification reports, etc.) emerged from the test case we chose to proceed. The last section concludes with suggested areas for future work and the main conclusions drawn from the work of the specific project.

2. Mission

Text classification, also known as text categorization, is a classical problem in natural language processing (NLP), which aims to assign labels or tags to textual units such as sentences, queries and documents. Machine learning models have drawn a lot of attention in recent years. Most classical machine learning based models follow the popular two-step procedure, where in the first step some hand-crafted features are extracted from the documents (or any other textual unit), and in the second step those features are fed to a classifier to make a prediction. Some of the popular hand-crafted features include bag of words (BoW), and their extensions.

Popular choices of classification algorithms include Naïve Bayes, support vector machines (SVM), hidden Markov model (HMM), gradient boosting trees, and random forests. The two-step approaches have several limitations. For example, reliance on the hand-crafted features requires tedious feature engineering and analysis to obtain a good performance. In addition, the strong dependence on domain knowledge for designing features makes the method difficult to easily generalize to new tasks. Finally, these models cannot take full advantage of large amounts of training data because the features (or feature templates) are pre-defined.

A paradigm shift started occurring in 2012, when a deep learning-based model, AlexNet¹ won the ImageNet competition by a large margin. Since then, deep learning models have been applied to a wide range of tasks in computer vision and NLP, improving the state-of-the-art². These models try to learn feature representations and perform classification (or regression), in an end-to-end fashion. They not only have the ability to uncover hidden patterns in data, but also are much more transferable from one application to another. Not surprisingly, these models are becoming the mainstream framework for various text classification and argumentation mining tasks in recent years. Argumentation mining aims to automatically identify structured argument data from unstructured natural language text. One particularly important aspect of argumentation mining is claim and evidence identification arising from abstracts of scientific articles, papers and documents, in order to extract useful information needed for better decision making and accurate planning controls.

Therefore, in this project, our main goal is to exploit unstructured parsing information to detect claims and evidences coming from citations for biomedical literature from MEDLINE, life science journals and online books and which are referring in keywords related with the targets and indicators of the specific SDG #3. Three (3) deep learning neural networks have been deployed as we mentioned in the previous section for the implementation of the specific approach and which have shown exceptional performances in various text classification tasks, including sentiment analysis, news categorization, topic classification, question answering (QA), and natural language inference (NLI), over the course of the past six years.

¹ <https://en.wikipedia.org/wiki/AlexNet>

² https://en.wikipedia.org/wiki/State_of_the_art

2.1 Previous efforts

There have been previous efforts to classify UN publications and facilitate document discovery and analytics. Specifically, DESA's Working Papers have recently been manually classified according to individual SDGs. There have also been a number of recent in-depth analyses of UN texts. Le Blanc, Freire, and Vierros (2017)³, for example, use a large collection of UN publications and academic sources to manually determine the connections among the ten targets of SDG 14. Vladimirova and Le Blanc (2015)⁴ used 40 global reports to carefully examine the links between education and other SDGs in flagship publications of the United Nations system. Le Blanc (2015)⁵ analyzed the targets in each of the 17 SDGs that refer to multiple goals and show the connections between some thematic areas. In each of these novel papers, the authors demonstrated the power of expert analysis and careful reading of individual texts to derive important insights.

However, there are limits to how well this methodology can scale and how it can be replicated with other texts. For any significant number of texts, the time and focus needed to understand them all becomes prohibitive. The problem gets worse as the number of documents continues to grow and as one discovers new connections between topics. For example, a publication that discusses inequality touches upon unemployment, gender, social protection, vulnerability, public policy, and many other relevant topics. Moreover, major publications like DESA's World Economic and Social Survey cover a broad range of topics related to development and simultaneously address multiple SDGs. As the Latin and Greek aphorism tells us, art is long, life is short.

Furthermore, additional efforts have been made also from Sovrano, Palmirani and Vitali (2004)⁶ as well as from Rodriguez Medina (2019)⁷. Specifically, as far as it concerns the first mentioned paper, they tried to produce a useful software for UN, that could help to speed up the process of qualifying the UN document following the SDGs in order to monitor the progresses at the world level to fight poverty, discrimination and climate changes using an automated labeling approach. As far as Rodriguez Medina (2019) is concerned, he tried to create and analyze a text classification dataset from freely-available web documents from the UN's SDGs and consequently used it to train and compare different multi-label text classifiers with the aim of exploring the alternatives for methods that facilitate the search of information of this type of documents.

³ Le Blanc, David, Clovis Freire, and Marjo Vierros. 2017. "Mapping the Linkages between Oceans and Other Sustainable Development Goals: A Preliminary Exploration." DESA Working Paper 149 (February). <https://www.un.org/development/desa/publications/working-paper/wp149>

⁴ Vladimirova, Katia, and David Le Blanc. 2015. "How Well Are the Links between Education and Other Sustainable development Goals Covered in UN Flagship Reports? A Contribution to the Study of the Science-Policy Interface on Education in the UN system." DESA Working Paper 146 (October). <https://www.un.org/development/desa/publications/working-paper/education-and-sdgs-in-un-flagship-reports>

⁵ Le Blanc, David. 2015. "Towards Integration at Last? The Sustainable Development goals as a Network of Targets." DESA Working Paper 141 (March). <https://www.un.org/development/desa/publications/working-paper/towards-integration-at-last>

⁶ Sovrano, Palminari, Vitali. 2004 "Deep Learning Based Multi-Label Text Classification of UNGA Resolutions". <https://arxiv.org/ftp/arxiv/papers/2004/2004.03455.pdf>

⁷ Rodriguez Medina, Samuel. 2019. "Multi-Label Text Classification with Transfer Learning for Policy Documents: The Case of the Sustainable Development Goals". <http://www.diva-portal.org/smash/record.jsf?pid=diva2%3A1360968&dswid=3657>

2.2 Approach of the problem

At this point, it seems logical as well necessary also, to present and elaborate further, our approach concerning the specific project. As we mentioned before, the main goal of this study, is to accurately predict the claims and evidences that displayed in the abstract part of documents, related with the targets and indicators of the SDG #3 of the UN in order to help governments and institutions to get better decisions and eliminate the time they need to read a paper. Therefore, we must first present very thoroughly and extensively the handling of this problem. First of all, in order to acquire all the necessary documents for our project, we implemented a web scraping approach by extracting the necessary data from the PubMed full-text archive of biomedical and life sciences journal literature of the U.S. National Institutes of Health's National Library. After that we used the certified online language processing service of "inception" that is part of the National Infrastructure for Language Resources and Technologies in Greece (CLARIN:EL), in order to properly annotate all the disposal documents that we had initially scrapped.

Subsequently, after all the necessary procedures required to format these documents in an appropriate type of data, we used as we mentioned before, three (3) deep neural network models. Neural network models have been demonstrated to be capable of achieving remarkable performance in sentence and document modeling. Convolutional neural network (CNN), Recurrent neural network (RNN) or Long Short-Term Memory Recurrent neural network (LSTM) and Feed Forward neural network (MLP) are three mainstream architectures for such modeling tasks, which adopt totally different ways of understanding natural languages. Specifically, a MLP neural network, is an artificial neural network wherein connections between the nodes do not form a cycle. As far as it concerns, a CNN model for text classification purposes is different from that of a neural network because it operates over a volume of inputs. Therefore, each layer tries to find a pattern or useful information of the data. Regarding also the RNN model is a class of artificial neural networks where connections between nodes form a directed graph along a temporal sequence. This allows it to exhibit temporal dynamic behavior.

In this work we tried to set the best hyperparameters and learning rates to each model separately in order to acquire the best possible results and in order for each model to be trained as efficiently as possible. Furthermore, in order to define the structure of the models we used appropriate embedding and dense layers, as well as the necessary activation functions in order to get the corresponding activation values of each neuron and to normalize the outputs of the networks. It must be stressed that we followed a specific approach in all of our three models. Specifically, we implemented a stratified shuffle split and a k fold cross-validation approach, in order to split the data in k datasets and in each iteration, to keep outside one dataset at a time, in order to train the models with the rest datasets. In order to evaluate the best model, we chose to keep the model with the least validation loss and subsequently with the greatest test accuracy. With this way, we managed to acquire better recalling results and comparing with a least biased method the capability of the models to cope with unknown data and therefore we did not rely on a random split and a model to train. Finally, in order to check the predictions of the best selected model, we used as a blind test case, a bunch of abstracts sentences derived from the PubMed as well and we tried to predict the label of each unique sentence.

3. Data

In this section of the project we will thoroughly present all the data we used for our analyses and for the deployment of our models, as well as all the pre-processing methods used and the main statistical results that derived after the necessary data transformation procedures.

3.1 Web scrapping

First of all, as we mentioned in the previous sub section of Chapter 2, we implemented a web scrapping approach, where by setting appropriately all the necessary queries for the selection of the right keywords that are referring to the SDG #3 we managed to acquire all the abstract texts from the PubMed full-text archive of biomedical and life sciences journal literature, that were related with the targets and indicators of the specific SDG. Web scraping, web harvesting, or web data extraction is a data scraping procedure, used for extracting data from websites. Web scraping software may access the World Wide Web directly using the Hypertext Transfer Protocol, or through a web browser. Specifically, for the implementation of the specific procedure we used the Beautiful Soup library which is a Python library for pulling data out of HTML and XML files. Subsequently, we extracted from all the disposal pages the urls' and afterwards we extracted and stored also, each abstract's text and title in a separate txt file every time. Finally, we saved each unique filename, url and title of the documents in a corresponding csv file, which afterwards they constituted our metadata files.

3.2 Annotation procedure using the CLARIN:EL

After the previous procedure mentioned, we ended up with a total of 150 txt documents, where with the help of the CLARIN:EL certified online language processing service of "inception", each member of our team, successfully annotated all these documents by giving each sentence the appropriate label that characterize them. Specifically, during this procedure we had to annotate properly the discourse topic of each abstract document, its argument i.e. the claim and evidence sentence or sentences as well as the research theme of each document's sentence. To be more specific, the topic sentence is what the discourse is about. A topic was usually being found, asking questions in general terms e.g., "What the text is actually about?", "What it reports?". We wanted actually to capture the aim, objective or the purpose of the paper. The topic ought to be indicative of the discourse and not to generic. Furthermore, we carefully noted to not annotate introductory expressions, the annotation span ideally should be a maximal nominal phase in a sentence, adjectives ought to be included in the beginning of an annotation span and possible punctuations at the end of an annotation span to not be included.

Regarding the annotation of the argument, we consisted a phrase with one or more claims that something was, or should being, the case and the evidence. The argument layer as we mentioned consisted from two labels. The claim and the evidence label. Generally, in an essay for example, argument consist of claims and premises. Premises can be supportive of attacking a claim. In our discourse, and in the range of an abstract we did not face these distinctions. Regarding also the claim label, we expected to find them in the end of the abstracts in the conclusion section, if the paper had sections, or in a sentence that concludes the results of the study. We should mentioned also that, some abstracts did not include claim sentences and for that reason we exclude them from our analysis. As far as the evidence is referred to, it was the reason for accepting the conclusion. Generally, the evidence may include experiments, observations, experiences or ideas, and may not be work that the authors have conducted themselves. The evidence in our corpus was usually seen as a result, and in many cases a result from a statistical test. The difficulty here, was to distinguish evidences from results. The rules that we followed were the following. Many abstracts report numerical data and results from statistical test. Sentences with raw numerical data that informed about the population for example, were not evidences. On the other hand, the statistical results (significant or not) were very likely evidences that support a conclusion. Furthermore, claim and evidence sentences ought to share keywords, ideally also found in topic, and/or the title of the paper.

Finally, regarding the research theme, we had to choose only one from a list of ten (10) disposal layers. Specifically, these were the following: Clinical Trial, Device, Diagnostic Tool, Drug, Infrastructure, Material, Prototype, Study, Treatment and Other. The Research Theme expressed the overarching goals of the authors and here we investigated how do they achieve it. Clinical trials are research studies performed in people that are aimed at evaluating a medical, surgical, or behavioral intervention. A medical device is any device intended to be used for medical purposes. Diagnosis is the identification of the nature and cause of a certain phenomenon. Infrastructure is defined as “the basic facilities, services and installations needed for the functioning of a community or society”. A drug is any substance that causes a change in an organism’s physiology or psychology when consumed. A biomaterial (or material in general) is any substance that has been engineered to inter-act with biological systems for a medical purpose - either a therapeutic or a diagnostic one. A prototype is an early sample, model, or release of a product built to test a concept or process. An observational study, investigators assess health outcomes in groups of participants according to a research plan or protocol. Finally, a therapy or medical treatment is the attempted remediation of a health problem, usually following a diagnosis.

After the whole annotation procedure was done, we implemented a curation procedure in order to correctly annotate any documents that had no agreement according to the majority of the annotators that were also annotated from another part of a different team their documents. This procedure was the final annotation procedure in order to conclude to documents with a powerful agreement. Therefore, after this procedure was finally ended, our final dataset was consisted from 900 csv files, where each file was a different abstract document and each row was consisted from the corresponding sentence of each document and from its characteristic label. More specifically, sentences that were lack of argument were characterized with “-”, and each evidence and claim sentence were characterized with the corresponding label, “claim” or “evidence”.

3.3 Data Parsing

After having our final dataset in our hands, i.e. the 900 different csv files, with the corresponding format that was mentioned before, we had to parse them accordingly in order to be able to analyze them in a properly manner. For the specific purpose we build a corresponding parser code, where we set appropriately each document to be placed in a different row into a unique data frame. Specifically, our data frame was consisted from two (2) columns. In the first column the label of the sentence was placed, while in the second column the sentence of each unique document was placed as well. Therefore, we ended up with 11064 rows and 2 variables. In the following table, we can get a quick view of how our data was originally seemed (Table 1).

Table 1: Structure of the data

Label	Sentence
NO LABEL	Multivariate Granger causality between CO2..
NO LABEL	Abstract:
NO LABEL	This paper addresses the impact of both..
EVIDENCE	The causality results indicate that there exis...
NO LABEL	The evidence seems to support the pollution
CLAIM	Overall, the method of managing both...
NO LABEL	Electricity consumption-GDP nexus in

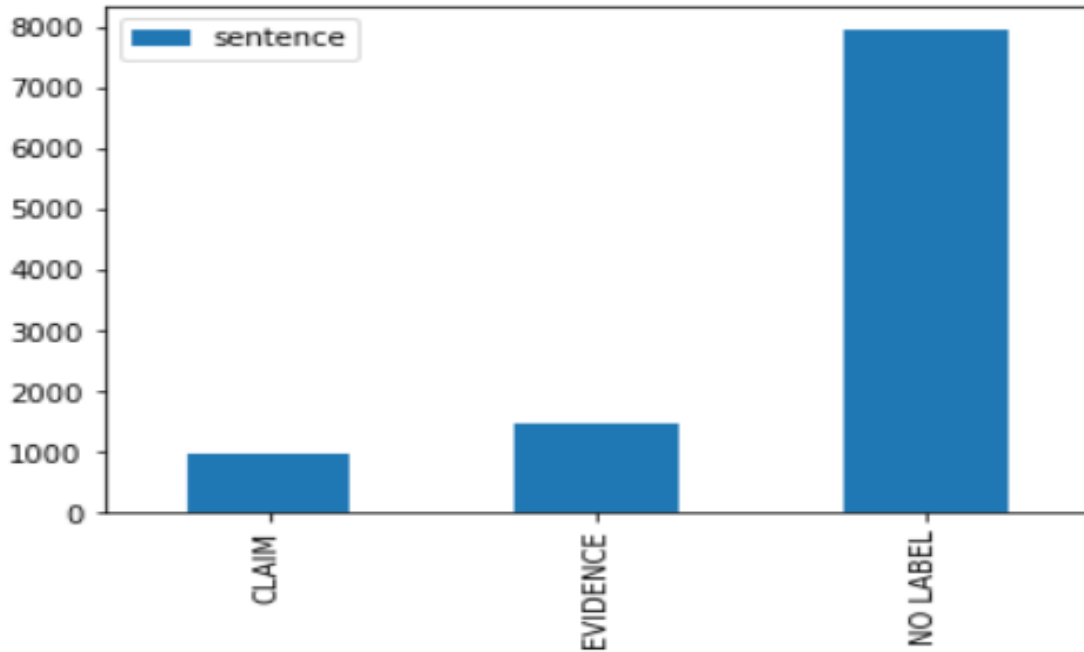
After having a quick view in our data, we assumed that they were very imbalanced and as a matter of fact, 77,62% (8588 sentences) of our data were classified with no label and only 13,43% (1486 sentences) and 8,95% (990 sentences) were classified as evidence and claim accordingly. After that, we took special care so that our data did not consist of missing values and as we noticed they did not exist at all. Furthermore, as we can observe in the following table (Table 2), there were many duplicate sentences and specifically the word 'Abstract' appeared 171 times inside our sentences. Therefore, we proceeded to their removal, as they would probably affect our final results.

Table 2: Main statistical results of the data

	Count	Unique	Top	Freq
LABEL	11064	3	NO LABEL	8588
SENTENCE	11064	10403	Abstract	171

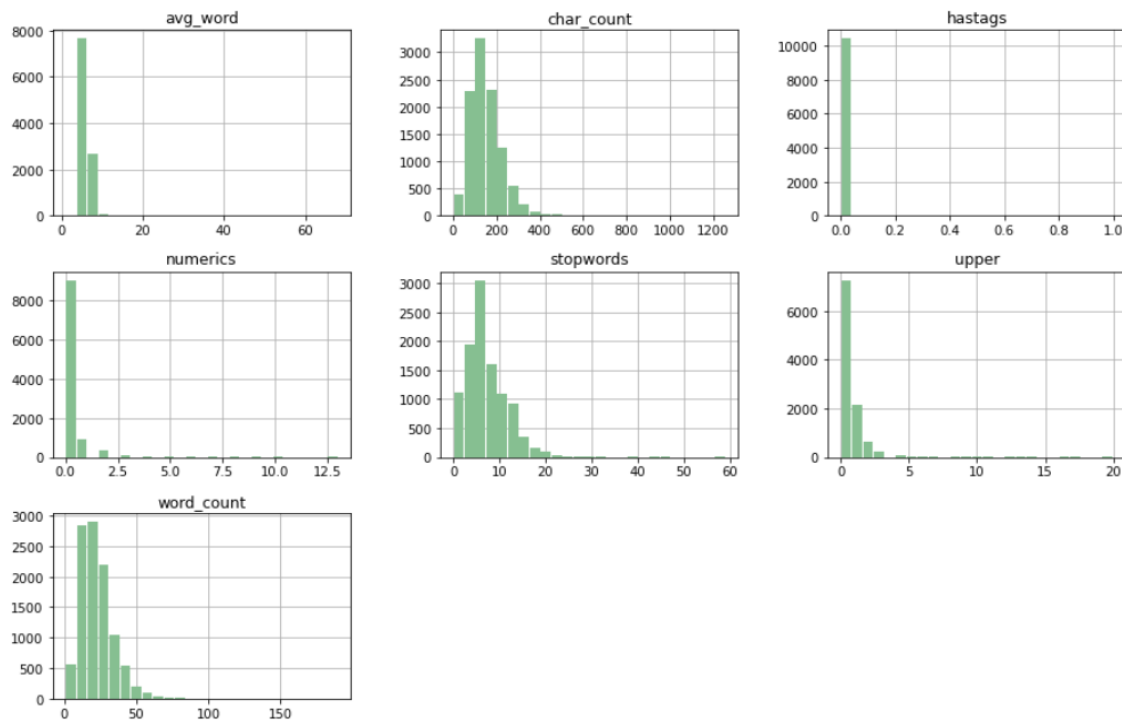
Subsequently, after the specific procedure, we ended up with 10403 observations. As far as the distribution of our data is concerned after the final removal of the duplicated sentences, as we can observe from the following bar plot (Plot 1), we can assume that 7953 observations were classified as 'No Label' sentences, 1470 as 'Evidence' sentences and 980 as 'Claim' sentences. In plain words, 635 sentences with the 'No label' layer were removed, 16 'evidence' sentences and 10 'claim' sentences were removed as well.

Plot 1: Distribution of the data



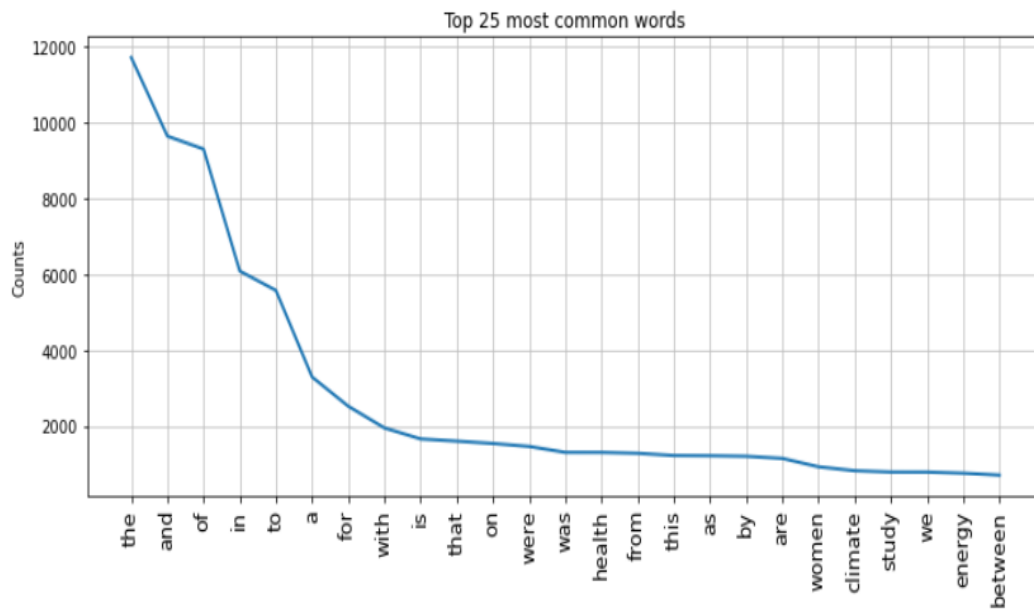
In order to have also a much more complete image of our dataset, we properly visualize using the following histograms (Plot 2), the distribution of the most common characteristics of our dataset, i.e. the number of words, characters, stop words, hashtags, numbers and upper cases letters that appeared in our sentences. As we can observe, the specific common characteristics, are following a right skewed distribution, which in plain words means that a large number of our data values occur on the left side with a fewer number of data values on the right side. A right-skewed distribution usually occurs when the data has a range boundary on the left-hand side of the histogram. Therefore, by taking a closer look to the specific plot (Plot 2), the majority of our sentences were consisted from 10 - 30 words and on average from 6 – 8 words per sentence. Furthermore, the number of characters and the number of stop words, were consisted in their majority, from 50 – 250 characters and 8 – 15 stop words long, accordingly. Finally, the number of hashtags, upper cases letters and numbers inside the sentences were very few, from the matter of fact as we can see, their majority were consisted under 5 unique corresponding abbreviations.

Plot 2: Distribution of the most common characteristics of the sentences



Continuing our analysis, we found also the most common words that appeared in our sentences. Specifically, as we can observe from the following line plot (Plot 3), the 25 most common words were the following, with the vast variety of them to be consisted from the words 'the', 'and', 'of', 'in' and 'to'.

Plot 3: 25 most common words



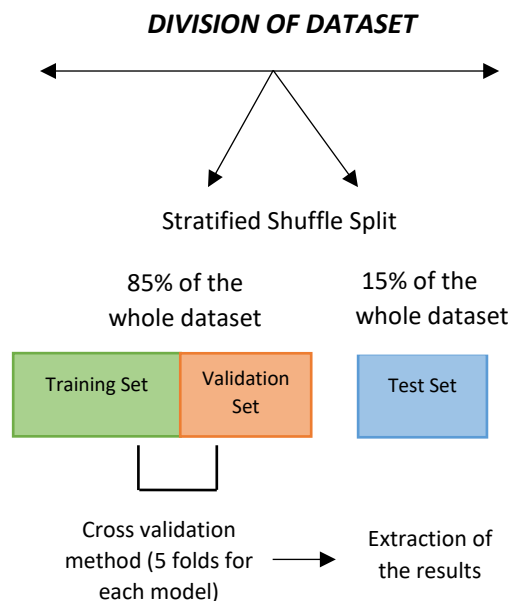
3.4 Pre-processing procedures

The first step that we had to implement before the deployment of our models, were any pre-processing procedures in the format of our data. Specifically, as a first step we proceeded to the conversion of all the upper cases letters to lower cases and we cleaned the text from any unwanted symbols, such as special characters, single characters and numbers that were inside the sentences. Furthermore, we removed all the predefined stop words, from the matter of fact, as we saw in a later stage, they affected our results. After the implementation of the specific procedures, we had to declare properly the input variable (X) and the predicted output variable (y). As it is logical, our input variable were all the sentences of the documents and the variable that we wanted to predict was obviously the corresponding label of the sentences.

For all the modelling procedures a specific intervention logic was used. Specifically, before deploying our models, we tried to split our initial dataset into two parts. The first part was divided by 85% of the initial dataset (train/validation dataset), which was used for the training and validation of our models, while the rest 15% was used for the predictive evaluation ability of the best model that was finally selected, which concerned data that had never been trained before. It should be mentioned also that the methodology we tried to follow for the experimental procedures of our models, was the stratified shuffle split and the cross-validation approach where we implemented 5 folds per iteration and the train – validations set were also split inside these folds. In statistics stratified shuffle sampling split is a sampling technique that is best used when a statistical population can easily be broken down into

distinctive sub-groups. Furthermore, the data is shuffled every time and then are also split. In addition, cross-validation is any of various similar model validation techniques for assessing how the results of a statistical analysis will generalize to an independent data set. It is mainly used in settings where the goal is prediction, and one wants to estimate how accurately a predictive model will perform in practice. Subsequently, in order to evaluate the best model, as we mentioned in a previous section of this report, we chose to keep the model with the least validation loss and subsequently with the greatest test accuracy. With this way, we managed to compare with a least biased method the capability of the models to cope with unknown data and therefore we did not rely on a random split and a model to train. Finally, in order to check the predictions of the best selected model, we used as a blind test case, a bunch of abstracts sentences derived from the PubMed library as well, and we tried to predict the label of each unique sentence.

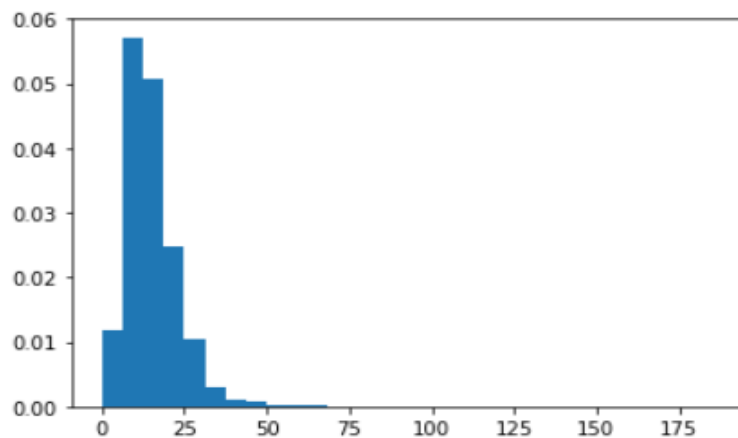
Plot 4: Division of training, validation and test datasets



It must be stressed that during the splitting procedures, we fitted as well the tokenizer on the train dataset only and we proceeded also to a padding procedure. At this point it seems appropriate to define what tokenization and padding is. Tokenization is a common task in Natural Language Processing (NLP). It's a fundamental step in both traditional NLP methods like Count Vectorizer and Advanced Deep Learning-based architectures. Given a character sequence and a defined document unit, tokenization is the task of chopping it up into pieces, called tokens, perhaps at the same time throwing away certain characters, such as punctuation. Therefore, in our project we used tokenization as a procedure of a sequence of characters in our documents which are grouped together as a useful semantic unit for the later on processing.

In addition, in any raw text data such as ours, naturally there will be sentences of different lengths. However, all neural networks require to have inputs with the same size. For this purpose, padding is done. Padding uses arguments such as sequences, padding, maxlen, truncating, value and dtype. In our case, we used together the 'post' padding and truncating procedure with the maximum length of the words to be specified as follows. We found that the maximal common length for the 99% of all sentences, were 43. This is also illustrated from the following histogram plot (Plot 5), where as we can observe, the maximum length of the words is lower than 70 words and specifically, under 65 words.

Plot 5: Distribution of the maximal length of the sentences



Therefore, for the implementation of our rest procedures we set as default the maximum length of the words to be 65 words long and as we observed the total number of words in the vocabulary was equal with 12048. Therefore, in our project by using the 'post' padding procedure, we set the padded values at the end of the sequences to be equal with 0 and by using also the 'post' truncating procedure as well, we truncated the sequences at the end. By use of maxlen to be equal with 65, we truncated the length of the padded sequences to 65. As a last step, before entering to the deployment of our models and consequently to the methodology that we followed, we used also a One-Hot-Encoding approach to the train data of our target variable, i.e. the Label variable, in order to convert them into a numerical data representation. Specifically, the integer encoded variable was removed and a new binary variable was added for each unique integer value inside the train dataset of our target variable.

4. Methodology

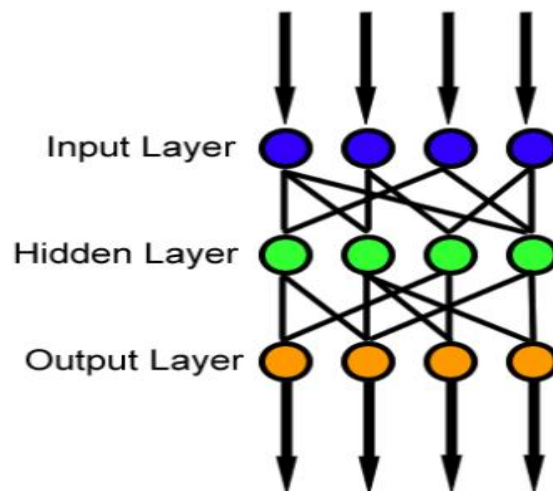
After having specified, the data we used for our analyses, as well as all the pre-processing procedures on the format of the data and the intervention logic behind the deployment of our models, we can now present all the methodologies specified for the architecture, structure, parameters, hyperparameters and optimizers of our models. Specifically, in this project, three different models will be created to successfully identify the type of each sentence as we have already mentioned. Each model has its specific architecture and hyperparameters. It is noteworthy that the training data set is imbalanced as it has already been stated in a previous section. Therefore, to train unbalanced classes 'fairly', we want to increase the importance of the two under-represented classes ("CLAIM", "EVIDENCE"). For the implementation of this procedure, we used a function from Keras, which automatically computes the class weights.

4.1 Models architectures and hyperparameters

1. Feed Forward Neural Network

The first model that was created, is a Feed Forward network. Feedforward Neural Network (FFNN) model is called a multilayer perceptron network that consists of input, hidden, and output layers of nonlinearly-activating nodes, which are interconnected in a feed-forward way with certain network weights. In FFNN, all the network weights are assigned random values initially, and the goal of the training is to adjust the set of network weights in a way that causes the output of the network to match the real values of the dataset as closely as possible (Figure 1).

Figure 1: Structure of a Feed Forward Neural Network



To increase the performance of the model in terms of accuracy, an embedding layer was inserted into the architecture of the model so that it could learn the relations between unique words of the input data. It should be noted that while using this embedding method, each word is mapped to one vector instead of just an integer, and the vector values are learned jointly with the neural network model. So, similar words in a semantic sense have a smaller (cosine) distance between them than words that have no semantic relationship. The embedding layer is used on the front end of a neural network and is fitted in a supervised way using the Backpropagation algorithm.

The hyperparameters that were changed and used for to create the input and embedding layer of the model were the following (Table 3):

Table 3: Hyperparameters of the Feed Forward Neural Network

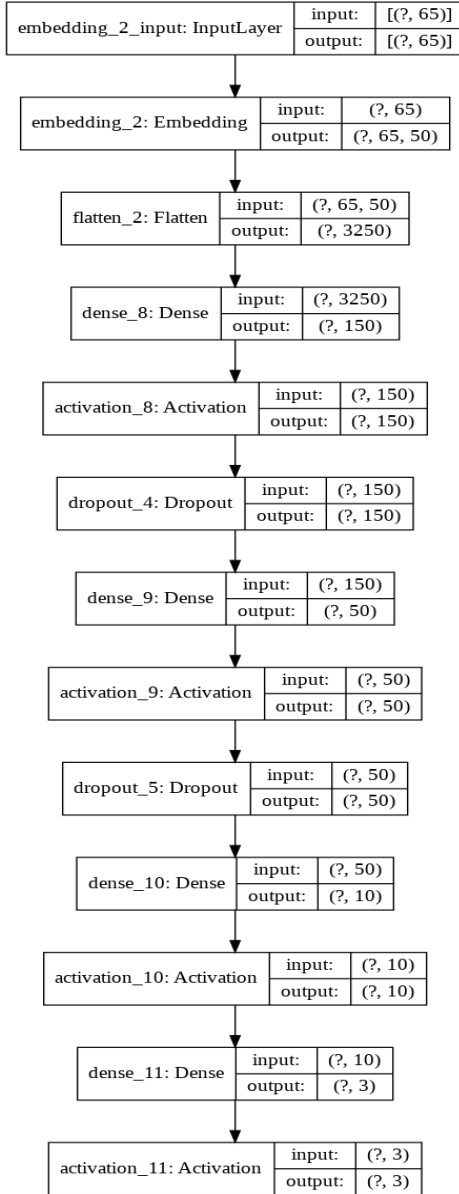
Name of hyperparameter	Value	Description of hyperparameter
<i>max_features</i>	12048	This number indicates the size of the vocabulary
<i>maxlen</i>	65	The length of each sentence
<i>embedding_dims</i>	50	The size of each word embedding
<i>batch_size</i>	64	The number of training sentences in one forward/backward pass of the model
<i>epoch number</i>	30	The number of complete passes through the training dataset.

During each propagation, 64 vectorized sentences of length 65 are inserted into the model. The embedding layer changes the shape of the input data from 1D to 2D, meaning that the shape of the samples will become (65x50) (50 represents the embedding dimensions of each word).

A. Model architecture

The following plot (Figure 2) visualizes the architecture of the FFNN model that was trained to classify the data. The specific model consists of an input layer that receives the vectorized data, whose length is equal to 65. Then the embedding layer converts the 1D shape of the data to a 2D shape, as it transforms each word into a vector of size 50. After the embedding layer, a Flatten layer is applied to flatten the input, meaning that the 2-dimensional input data will be transformed to 1D, without losing the information that was added by the embedding layer (the semantic relationship between the words). Next, we add 3 dense layers and between them drop-out layers. The first dense layer consists of 150 neurons, the second layer consists of 50 neurons, and the third 10 neurons.

Figure 2: Architecture of feed-forward neural network



The 3 dropout layers that are located after each dense layer, have a drop-out rate of 30%. As a result, during training, some number of layer outputs are randomly ignored or “dropped out” with a probability of 30%. This has the effect of making each dense layer be treated as a layer with a different number of nodes. This indicates that the neural network is forced to find new ways to classify the data, thus reducing the chance of overfitting. Without any dropout, our neural network exhibits substantial overfitting.

Lastly, the output layer is added and is the one that outputs the classification predictions for each sentence. To be more specific it outputs a vector with a length of three since this project aims to classify the sentence into one of three classes: “NO LABEL”, “CLAIM”, “EVIDENCE”. Each element of this vector is a probability, which shows the possibility of a sentence to belong to a specific class. It should also be noted that for each dense layer, an activation layer was applied with the Relu function. Since the problem of this project is a multiclass classification problem, we use the Softmax function to output the predictions.

B. Number of parameters

The total number of parameters, which are tuned during the learning procedure of the model, is also calculated in the specific paragraph. Since each distinct word is represented as a 50-sized vector, the total number of parameters that are created by the embedding layer is equal to:

$$\text{vocabulary size} \times \text{embedding dimensions} = 12048 \times 50 = 602400$$

The input data that the Flatten layer receives has shape (none, 65, 50) and then it is converted to (none, 65x50 = 3250). Each value of the resulting input data is connected with each neuron of the 1st hidden layer via a specific weight. Therefore, the number of parameters that are going to be created in the 1st layer is:

$$\text{input dimensions} \times \text{number of neurons of 1st layer} = 3250 \times 150 = 487500$$

It is important to note that in each multiplication between the input and the weights of each neuron, another parameter is added to the equation, which is named “bias parameter”. Therefore, the total number of parameters in the 1st layer is the following:

$$\begin{aligned} &\text{input dimensions} \times \text{number of neurons of 1st layer} + \text{bias parameters} \\ &= 3250 \times 150 + 1 \times 150 = 487650 \end{aligned}$$

The input data that the 2nd hidden layer receives has shape (none, 300), where 300 are the dimensions of the 1st layer’s output. Therefore, the number of new parameters that have been added is:

$$\begin{aligned} &\text{input dimensions} \times \text{number of neurons of 2nd layer} + \text{bias parameters} \\ &= 150 \times 50 + 1 \times 50 = 7550 \end{aligned}$$

Furthermore, new parameters have been added from the 3rd layer and their count is:

$$\begin{aligned} &\text{input dimensions} \times \text{number of neurons of 3rd layer} + \text{bias parameters} \\ &= 50 \times 10 + 1 \times 10 = 510 \end{aligned}$$

Lastly, the total number of parameters increases by 33, due to the multiplication of the output dimensions of the last hidden layer (none, 10) with the 3 neurons of the output layer. In total, the model contains 1098143 trainable parameters.

C. Optimizer and Evaluation Metrics

For the optimizer, we used the Adam optimization algorithm. Since this model is sensitive to overfitting, we also decreased the learning rate to 0,005%. Note that the learning rate controls how quickly or slowly the neural network model learns the data. A smaller learning rate may allow the model to learn a more optimal or even globally optimal set of weights but may take significantly longer to train. For evaluation metrics we made use of the Categorical Cross-Entropy loss function and categorical accuracy.

2. Convolutional Neural Network

The second model that was created, was a convolutional network with an embedding layer. In a traditional feedforward neural network, each input neuron is connected to each output neuron in the next layer. In contrast, the CNN model uses convolutions over the input layer to compute the output. This convolutional network contains one 1D convolutional layer. This convolutional layer contains a series of filters known as convolutional kernels. Each convolutional kernel is like a sliding window, whose job is to look at embeddings for multiple words. Each kernel is designed to look at a word, and surrounding words in a sequential window, and output a value that captures something about that phrase. To set up a network so that it is capable of learning a variety of different relationships between words, it is crucial to provide the model with many filters (Figure 3). The hyperparameters that were adjusted to create the input, embedding layer, and the 1D convolutional layer of the model are the following (Table 4):

Figure 3: Structure of a Convolutional Neural Network for NLP

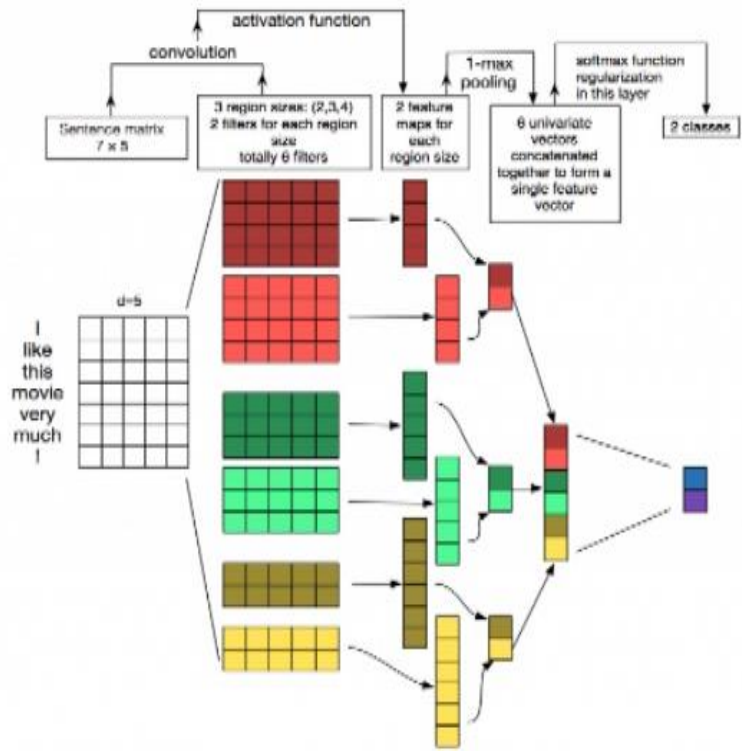


Table 4: Hyperparameters of the Convolutional Neural Network

Name of hyperparameter	Value	Description of hyperparameter
max_features	12048	This number indicates the size of the vocabulary
maxlen	65	The length of each sentence
embedding_dims	50	The size of each word embedding
nof_filters	64	The number of kernels
kernel size	3	The height of each kernel

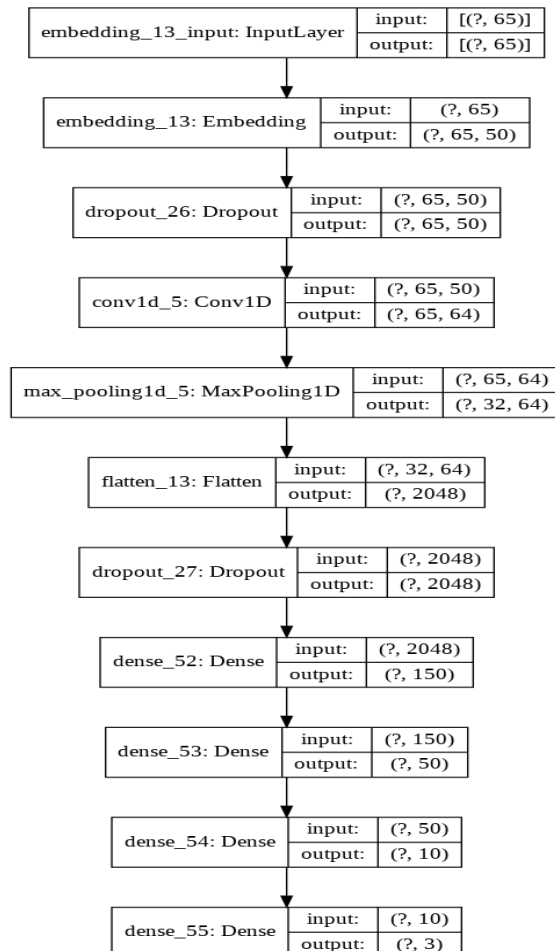
<i>batch_size</i>	64	<i>The number of training sentences in one forward/backward pass of the model</i>
<i>epoch number</i>	30	<i>The number of complete passes through the training dataset.</i>

Note that the kernels will no longer be square, instead, they will be a wide rectangle with dimensions $\text{kernel_size} \times \text{embedding_dims} = 3 \times 50$. The height of the kernel will be the number of embeddings it will see at once, similar to representing an n-gram in a word model. In this case, it will look at each 3-gram in each sentence.

A. MODEL ARCHITECTURE

The following plot (Figure 4) visualizes the architecture of the CNN model that was trained to classify the data. The specific model consists of an input layer that receives the vectorized data, whose length is equal to 65. Then the embedding layer converts the 1D shape of the data to a 2D

Figure 4: Architecture of Convolutional neural network



shape, as it transforms each word into a vector of size 50. After the creation of the embedding layer, we set a Dropout layer with a dropout rate = 40%.

The output (sentence) of the embedding layer is shaped as a 65x50 matrix. This matrix acts as the input layer for the 1D convolutional layer. The elements involved in carrying out the convolution operation are the kernels/filters. Every filter performs convolution on the sentence matrix and generates (variable-length) feature maps. Since each kernel outputs a specific vector after the convolution, the width of the final output will be equal to the number of filters (64). It is worth mentioning that the output's height of the convolutional layer has not changed as it would be, if we didn't use padding. This means that if we didn't use padding, the output's shape of the convolutional layer would be (63, 64), meaning that the height of the final output would be 63 instead of 65. This happens because each filter is applied systematically to the input sentence matrix. It starts from the top of the matrix and is moved from top to bottom by one row at a

time until the edge of the filter reaches the bottom edge of the sentence matrix. For a 3x50 filter applied to a 65x50 sentence matrix, we can see that it can only be applied 63 times. The problem with this phenomenon is that the kernels didn't pay attention to the edges of the input matrix and disregarded that information (first and last word). The reduction in the size of the output to the feature map is referred to as *border effects*⁸. After the convolutional layer was applied, a 1-max pooling is performed over each map and the largest number from each feature map is recorded. According to the Keras Documentation, the MaxPooling1D layer⁹ "downsamples the input representation by taking the maximum value over the window defined by pool_size". In our case, the pool size is equal to two. So, the height of the output would be equal to the height of the input divided by two.

A Flatten layer is then used to flatten the input, resulting to an output that has a shape of (none, 32 x 64) = (none, 2048). A dropout layer is also applied with the same dropout rate as the previous one. The last four layers that we use are dense layers. The first layer contains 150 neurons, the second 50 neurons, and the third one 10 neurons. The last layer is the output layer which in turn outputs the classification predictions for each sentence.

B. Number of parameters

The total number of parameters, which are tuned during the learning procedure of the model, is also calculated in the specific paragraph. Since each distinct word is represented as a 50-sized vector, the total number of parameters that are created by the embedding layer is equal to:

$$\text{vocabulary size} \times \text{embedding dimensions} = 12048 \times 50 = 602400$$

To calculate the learnable parameters in the 1D convolutional layer, we have to multiply by the shape of width, height of current filters and number of previous layer's filters. There is also a bias term for each of the filter. In this case, there are no previous filters so the variable previous layer's filters is equal with zero.

Number of parameters in this CONV layer would be:

$$((\text{shape of width of the filter} \times \text{shape of height of the filter} + 1) \times \text{number of filters}) = ((3 \times 50 + 1) \times 64 = 9664$$

⁸ Jason Brownlee. A Gentle Introduction to Padding and Stride for Convolutional Neural Networks. machinelearningmastery.com. [Online] April 19, 2019 . <https://machinelearningmastery.com/padding-and-stride-for-convolutional-neural-networks/>.

⁹ MaxPooling1D layer. [Online] https://keras.io/api/layers/pooling_layers/max_pooling1d/.

After the Max-Pooling layer, the shape of the input will be (32, 64). The Flatten layer then is applied and the input's dimensions are equal to $32 \times 64 = 2048$. Lastly, 3 dense layers and the output layer are applied, which in turn add $307350 + 7550 + 510 + 33 = 315443$ new trainable parameters. In total, the model contains 927507 trainable parameters.

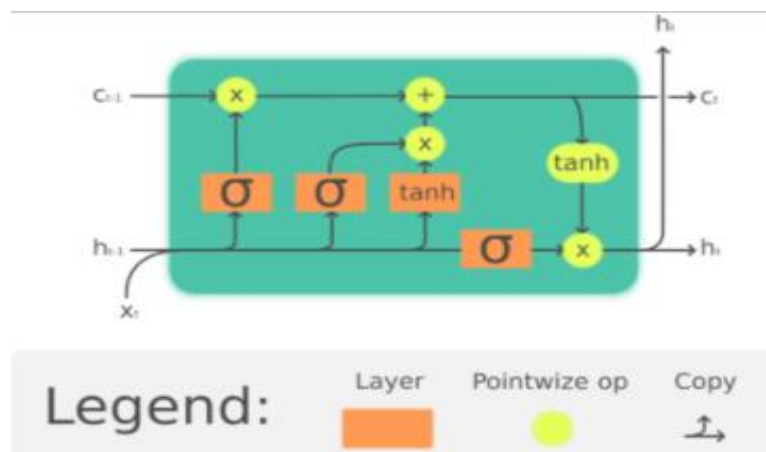
C. OPTIMIZER AND EVALUATION METRICS

For the optimizer, we used the Adam optimization algorithm. Since this model is sensitive to overfitting, we also decreased the learning rate to 0,005%. For evaluation metrics we made use of the Categorical Cross-Entropy loss function and categorical accuracy.

3. Long Short-Term Memory Neural Network

The third model that was created, is a Bidirectional LSTM model. A Long short-term memory (LSTM) is an artificial recurrent neural network (RNN) architecture, which according to Wikipedia¹⁰, manages to “deal with the vanishing gradient problem that can be encountered when training traditional RNNs”. Traditional RNNs suffer from the problem of vanishing gradients, which hampers learning of long data sequences. The advantage of an LSTM cell compared to a common recurrent unit is its cell memory unit. The cell vector has the ability to encapsulate the notion of forgetting part of its previously stored memory, as well as to add part of the new information. Bidirectional LSTMs are an extension of traditional LSTMs that can improve model performance on sequence classification problems. Bidirectional LSTMs train two instead of one LSTMs on the input sequence. The first on the input sequence as-is and the second on a reversed copy of the input sequence. This can provide additional context to the network and result in faster and even fuller learning on the problem (Figure 5).

Figure 5: Structure of a Long Short-Term Memory Neural Network



¹⁰ Long short-term memory. Wikipedia. [Online] https://en.wikipedia.org/wiki/Long_short-term_memory.

The hyperparameters that were adjusted to create the architecture the RNN model are the following (Table 5):

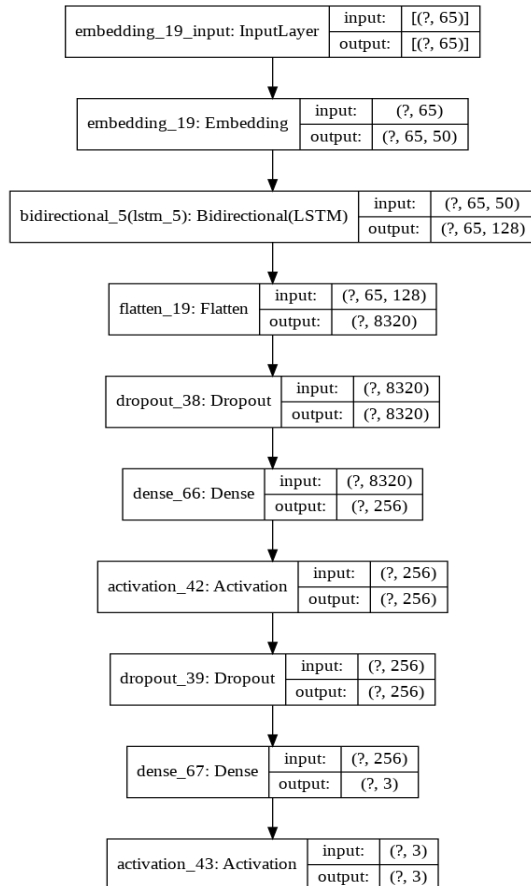
Table 5: Hyperparameters of the Long Short-Term Memory Neural Network

Name of hyperparameter	Value	Description of hyperparameter
max_features	12048	This number indicates the size of the vocabulary
maxlen	65	The length of each sentence
embedding_dims	50	The size of each word embedding
batch_size	32	The number of training sentences in one forward/backward pass of the model
epoch number	10	The number of complete passes through the training dataset.
memory units	64	The number of memory units for each LSTM hidden cell

It should be mentioned that one bidirectional LSTM layer will be created, meaning that two instead of one LSTMs will be trained on the input sequence and each one will contain hidden cells, which will in turn contain 64 memory units (in total 128 units).

A. Model Architecture

Figure 6: Architecture of Bidirectional LSTM model



The specific model (Figure 6) consists of an input layer that receives the vectorized data, whose length is equal to 65. Then the embedding layer converts the 1D shape of the tensor to a 2D shape, as it transforms each word into a vector of size 50. In Bi-LSTM there will be one LSTM unrolling from left to right (LSTM1) on the input sentence matrix and another LSTM unrolling from right to left (LSTM2). Since the input size is $(n, 65, 50)$, where n refers to the batch size, LSTM1 will return output of size $(n, 65, 64)$ and LSTM2 will also return output of size $(n, 65, 64)$.

The layer will then combine the output of LSTM1 and LSTM2 by applying element wise concatenation of LSTM1 output to LSTM2 at each timestep. This will result in an output of size $(n, 65, 128)$. A Flatten layer is then used to flatten the input, resulting to an output that has a shape of $(\text{none}, 65 \times 128) = (\text{none}, 8320)$. After the creation of the Flatten layer, we set a Dropout layer with a dropout rate = 50%. A dense layer is also added to the model and it contains 256

neurons. A dropout layer is also applied with the same dropout rate as the previous one. The last layer is the output layer, which outputs the classification predictions for each sentence.

B. NUMBER OF PARAMETERS

The total number of trainable parameters that are created by the embedding layer is equal to:

$$\text{vocabulary size} \times \text{embedding dimensions} = 12048 \times 50 = 602400$$

The number of trainable parameters that are created by the bidirectional LSTM layer are calculated by the following equation:

$$4 \times ((\text{size of input} + \text{bias term}) \times \text{size of output} + \text{size of output}^2) \times 2$$

The $\times 2$ is added because the bidirectional LSTM contains two LSTMs.

Therefore, the number of new trainable parameters is equal to:

$$4 \times ((50 + 1) \times 64 + 64^2) \times 2 = 58880$$

The Flatten layer that is applied later, returns a vector output of length 6400 for each sentence. The dense layer which is located after the Flatten layer has 256 neurons. This means that the number of new trainable parameters from the Dense layer will be equal to:

$$8320 \times 256 + 1 \times 256 = 2130176$$

The last new trainable parameters will be created in the output layer and will number:

$$256 \times 3 + 1 \times 3 = 771$$

In total, the model contains 2792227 trainable parameters.

C. Optimizer and Evaluation Metrics

For the optimizer, we used the Adam optimization algorithm. Since this model is sensitive to overfitting, we also decreased the learning rate to 0,005%. Note that the learning rate controls how quickly or slowly the neural network model learns the data. A smaller learning rate may allow the model to learn a more optimal or even globally optimal set of weights but may take significantly longer to train. For evaluation metrics we made use of the Categorical Cross-Entropy loss function and categorical accuracy.

4.2 *Model comparison*

The dataset as mentioned in a previous section, has been split into two sets. The train-validation dataset and the test dataset. The ModelCheckpoint callback function¹¹ is used in conjunction with the training of each model to save a model or weights (in a checkpoint file) at some interval, so the model (or weights) can be loaded later to test its performance. This function stores the best model where the validation loss is least among all epochs. We will use the test dataset, to test the performance of each model that was created. It is important to note that we will mostly concentrate on the accuracy of the minority classes to identify the best model, as well as from our business perspective to the model with the best recall evaluation results in the “CLAIM” and “EVIDENCE” labels. The final results of our models are thoroughly described in the next chapter.

¹¹ ModelCheckpoint. Keras. [Online] https://keras.io/api/callbacks/model_checkpoint/.

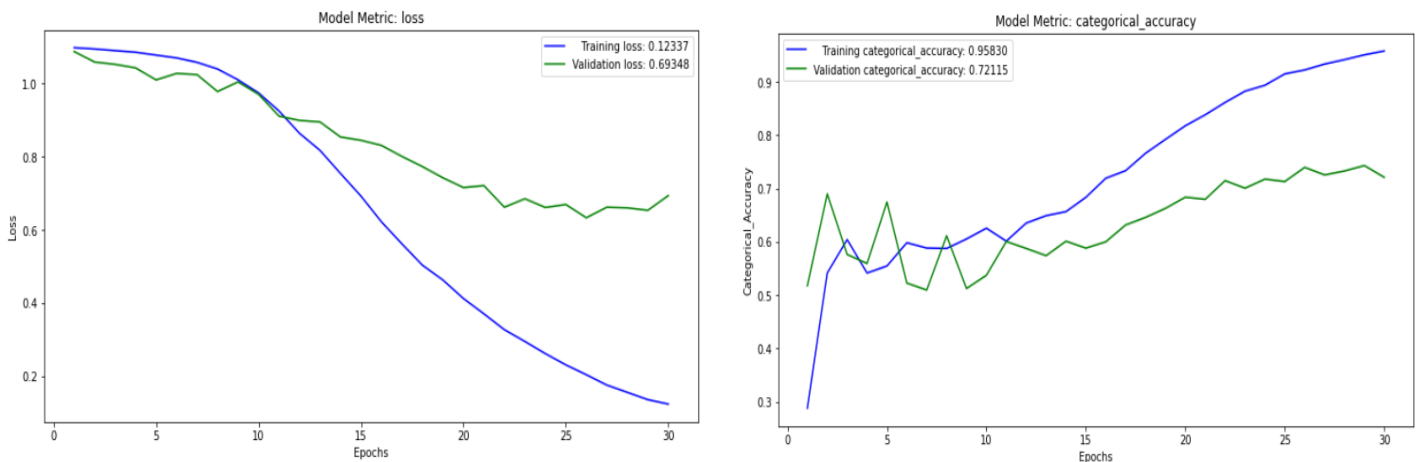
5. Results

5.1 Main findings

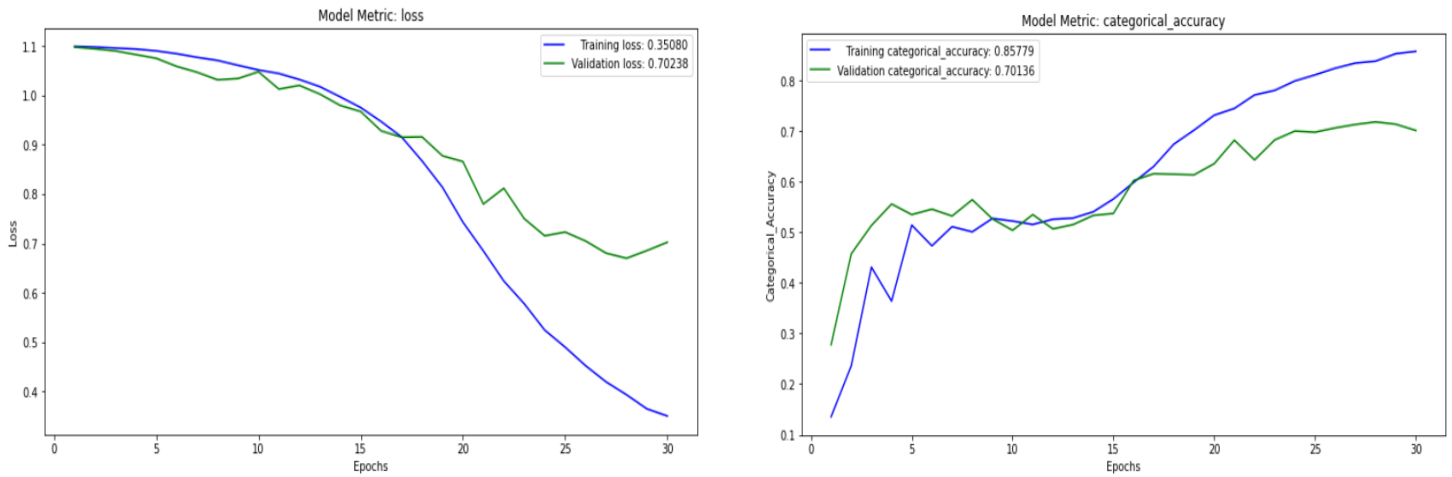
After having discussed and analyzed the architectures, the number of parameters and the optimizer and evaluation metrics used in all of our neural network models, we can at this point present all of our findings. Specifically, by plotting first the corresponding learning curves (Plots 6 – 8), that are referring to the training and validation losses as well as to the training and validation categorical accuracies, we can get a first look about the performance of our models. At this point we should mention, that a learning curve is a plot of model learning performance over experience or time and are a widely used diagnostic tool in machine learning for algorithms that learn from a training dataset incrementally. Reviewing learning curves of models during training can be used to diagnose problems with learning, such as an underfit or overfit model, as well as whether the training and validation datasets are suitably representative.

To be more precise, our main interest regarding the training and validation loss is when the curve of the validation loss starts to get increase after having reached to a local minimum. At this point we therefore select the number of the disposal epochs of the model, before starting overfitting. In statistics, overfitting is "the production of an analysis that corresponds too closely or exactly to a particular set of data, and may therefore fail to fit additional data or predict future observations reliably". The essence of overfitting is to have unknowingly extracted some of the residual variation (i.e. the noise) as if that variation represented underlying model structure. On the contrary, the training and validation categorical accuracy are like a mirror to the previous mentioned curves. Specifically, when the validation's categorical accuracy curve reaches to a local maximal, at this point we suppose that the number of the selected epochs is the number of the optimal epochs that the model should be use to fit and therefore to stop it from overfitting.

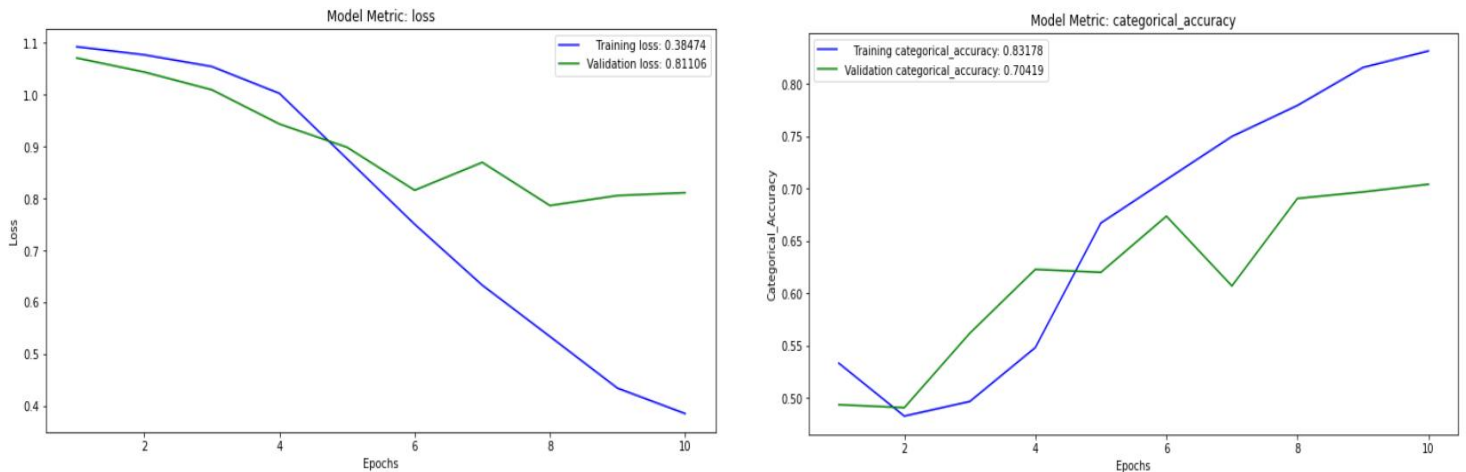
Plot 6: Learning curves of the MLP model



Plot 7: Learning curves of the CNN model

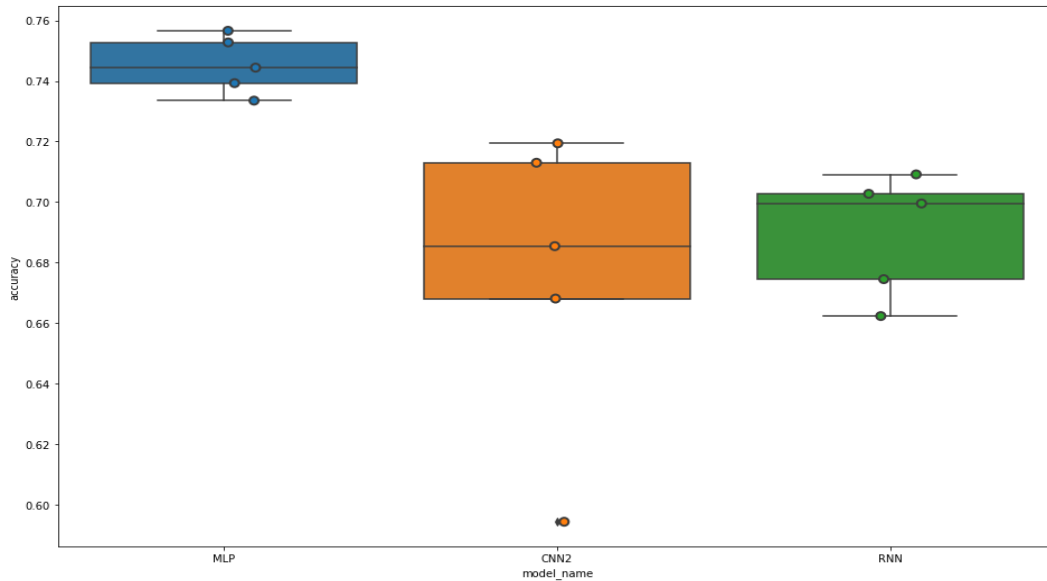


Plot 8: Learning curves of the RNN model



From what we can see from the above learning curves of our models, the MLP and CNN model present the local optima in the 30 epochs, while the RNN model during the 10 epochs. The greatest validation loss seems to be presented with the RNN model at 0.81, while the greatest validation categorical accuracy seems to be presented with the MLP model at 0.72. Continuing our analysis, and after having trained our models with 5-fold cross validation as we have already mentioned and by keeping the accuracy of each training, we ended up with the following box plots (Plot 9) containing the results of each model in the test data-set. From what we can observe, the plot below, indicates that the RNN and CNN models, are not the appropriate models to use for our business purpose and therefore it seems that the MLP model, is the model with the best results. This is obvious from the fact that, the MLP model has less variance in its predictions in comparison with the other two models which present a greatest variance to their predictions and especially the CNN model. Furthermore, another aspect to consider that the MLP model has better testing accuracies, is that its upper and lower whiskers are very close to the upper and lower quartiles accordingly and that the predictions are all very close to the median.

Plot 9: Test accuracies of all models using the cross-validation procedure



At this point we should recall that, in order to evaluate the best model, we chose to keep the model with the least validation loss and subsequently with the greatest test accuracy. With this way, we managed to acquire better recalling results and comparing with a least biased method the capability of the models to cope with unknown data and therefore we did not rely on a random split and a model to train. Although, accuracy from its' own, is not always the best metric for evaluating and selecting a model that makes good predictions for each label. In order to further analyze the competence of each model, we used several other metrics, mainly arisen from corresponding confusion matrices and classification reports. Specifically, except from the

accuracy of the model, we selected also to evaluate the predictive ability of our models, using the metrics of Precision and Recall (Figure 7). The first one (Precision), is referring to the percentage of the right predictions classified by each model for every label and it is calculated as follows.

$$\frac{TP}{TP + FP}$$

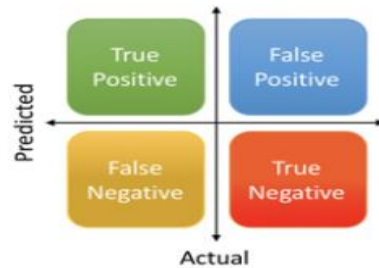
The other one (Recall), is referring to the percentage of the actual labels being predicted correctly (the sensitivity of our model in finding the evidence and claims) and it is calculated as follows.

Figure 7: Confusion Matrix plot and selected metrics

$$\text{Precision} = \frac{\text{True Positive}}{\text{Actual Results}} \quad \text{or} \quad \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

$$\text{Recall} = \frac{\text{True Positive}}{\text{Predicted Results}} \quad \text{or} \quad \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

$$\text{Accuracy} = \frac{\text{True Positive} + \text{True Negative}}{\text{Total}}$$



$$\frac{TP}{TP + FN}$$

Where,

TP = the number of true positives for a specific label

FP = the number of false positives for a specific label

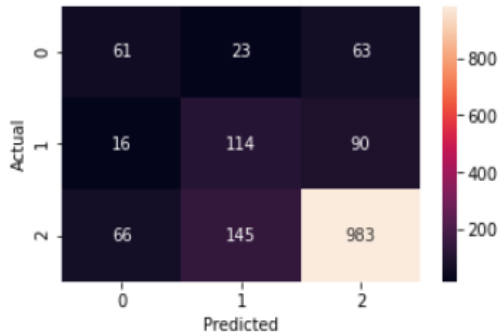
FN = the number of false negatives for a specific label

Both are important metrics because we wanted to choose a model that makes good predictions for each label. For our business purpose we selected to proceed with a model which predicts just as well, the number of

$$\frac{TP}{TP + FP} \quad \text{and} \quad \frac{TP}{TP + FN}$$

giving a greatest importance to the precision metric. Therefore, our final results were the following (Plots 10 - 12 and Tables 6 - 8).

Plot 10: Confusion matrix of the MLP model

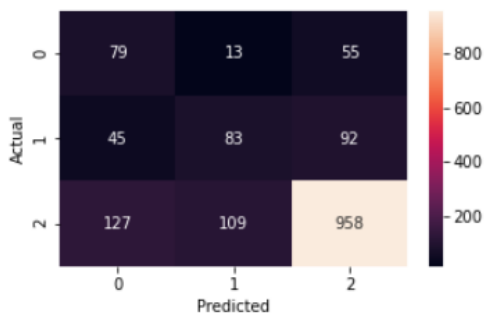


Test accuracy of MLP: 75.65%

Table 6: Classification report of the MLP model

	Precision	Recall	F1-score	Support
CLAIM	0.43	0.41	0.42	147
EVIDENCE	0.40	0.52	0.45	220
NO LABEL	0.87	0.82	0.84	1194
Accuracy			0.74	1561
Macro avg	0.57	0.59	0.57	1561
Weighted avg	0.76	0.74	0.75	1561

Plot 11: Confusion matrix of the CNN model

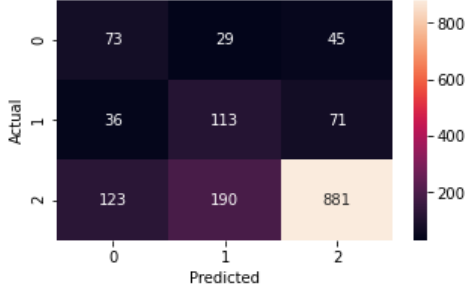


Test accuracy of CNN: 71.3%

Table 7: Classification report of the CNN model

	Precision	Recall	F1-score	Support
CLAIM	0.31	0.54	0.40	147
EVIDENCE	0.40	0.38	0.39	220
NO LABEL	0.87	0.80	0.83	1194
Accuracy			0.72	1561
Macro avg	0.53	0.57	0.54	1561
Weighted avg	0.75	0.72	0.73	1561

Plot 12: Confusion matrix of the LSTM model



Test accuracy of LSTM: 69.95%

Table 8: Classification report of the LSTM model

	Precision	Recall	F1-score	Support
CLAIM	0.31	0.50	0.39	147
EVIDENCE	0.34	0.51	0.41	220
NO LABEL	0.88	0.74	0.80	1194
Accuracy			0.68	1561
Macro avg	0.51	0.58	0.53	1561
Weighted avg	0.75	0.68	0.71	1561

In conclusion, MLP has 43% precision for claims and 40% for evidence, results better than those of CNN (with 31% for claims and 40% for evidence) and RNN (with 31% for claims and 34% for evidence). With respect to recall RNN has astoundingly better results with 50% and 51%. In this case if we wanted to have a model that misses less evidence and claims, but makes more misclassifications in the 'no label' then the RNN would be our choice and this is the reason that accuracy is not always the best metric to use. In our case for our business purpose, as we said before, precision is more important and for this reason MLP is the final model we chose, which has overall better precision and recall.

5.2 Blind test set predictions

After having selected the best model for our business purpose, we proceed to the labels' predictions, of the unseen data so far and which have obtained as we have also mentioned in previous section from the PubMed source. Specifically, after having correctly parsed again this new dataset, we proceed once again to the pre mentioned cleaning procedures as before and to the tokenization and padding (using post padding and post truncating) of the new sentences. After that, we load accordingly our final best selected model (MLP: 'final_mlp_model.h5') and make the corresponding predictions after having also proceeded to the encoding of the new target variable using a One Hot Encoder. As a last step we saved our predictions to a new column inside the blind test dataset as we can see below (Table 9).

Table 9: Blind test set predictions according to the best model

	sentence	file_origin_name	label_predictions
0	exploration corrective measures greenhouse gas...	ABC_G1B1_10.1016.j.jclepro.2019.118645.csv	NO LABEL
1	abstract	ABC_G1B1_10.1016.j.jclepro.2019.118645.csv	NO LABEL
2	greenhouse gas emission increasing alarmingly ...	ABC_G1B1_10.1016.j.jclepro.2019.118645.csv	NO LABEL
3	researchers struggling minimize emission using...	ABC_G1B1_10.1016.j.jclepro.2019.118645.csv	NO LABEL
4	paper presents mathematical model parameters r...	ABC_G1B1_10.1016.j.jclepro.2019.118645.csv	NO LABEL
5	greenhouse gas emission existing fossil fuel p...	ABC_G1B1_10.1016.j.jclepro.2019.118645.csv	NO LABEL
6	result shows power plants coal diesel natural ...	ABC_G1B1_10.1016.j.jclepro.2019.118645.csv	EVIDENCE
7	furthermore several greenhouse gas mitigating ...	ABC_G1B1_10.1016.j.jclepro.2019.118645.csv	EVIDENCE
8	exploration household energy choice expenditur...	ABC_G1B1_10.1016_j.energy.2017.06.117.csv	NO LABEL
9	globally 1 3 billion people access electricity...	ABC_G1B1_10.1016_j.energy.2017.06.117.csv	NO LABEL
10	projected 2030 1 billion people access electri...	ABC_G1B1_10.1016_j.energy.2017.06.117.csv	NO LABEL
11	widely documented use firewood biomass househo...	ABC_G1B1_10.1016_j.energy.2017.06.117.csv	NO LABEL
12	massive use firewood biomass lead deforestatio...	ABC_G1B1_10.1016_j.energy.2017.06.117.csv	NO LABEL
13	therefore imperative ensure use clean energy e...	ABC_G1B1_10.1016_j.energy.2017.06.117.csv	CLAIM
14	encourage use clean energy factors affect hous...	ABC_G1B1_10.1016_j.energy.2017.06.117.csv	NO LABEL
15	using information collected 29 000 households ...	ABC_G1B1_10.1016_j.energy.2017.06.117.csv	NO LABEL
16	particularly households use clean energy progr...	ABC_G1B1_10.1016_j.energy.2017.06.117.csv	NO LABEL
17	based findings study suggests specific policie...	ABC_G1B1_10.1016_j.energy.2017.06.117.csv	CLAIM
18	cooking season risk factors acute lower respir...	ABC_G1B1_10.1371_journal.pone.0128933.csv	NO LABEL
19	background	ABC_G1B1_10.1371_journal.pone.0128933.csv	NO LABEL

6. Conclusions and future work

Since the aim of this project was to identify the arguments of each abstract, we had placed great emphasis on the correct classification of the minority classes, hence we used the precision and recall metrics to evaluate the performance of the models, as well as select the best one. All the models that we trained, managed to capture a high total accuracy, but had low recall and precision metrics, meaning that they couldn't classify properly the minority classes. This could be attributed to the small training dataset that was produced, as well as, to the different methods each group of annotators used to manually identify the claims and evidences of each abstract. This could lead to models that overfit quickly. The smaller the dataset, the more difficult it is for each model to learn the problem. This issue could be solved, by using various regularization methods that modify our learning algorithm to reduce its generalization error but not its training error. Some common regularization methods that we used to somehow limit the problem of overfitting, were the insertion of dropout layers, the implementation of early stopping using callbacks, and the re-adjustment of class weights to better understand the performance of each model.

In our future plans, another regularization method that we could have used to limit overfitting, is Data Augmentation. The Data Augmentation technique could be applied, which increases the size of the training data, by applying random (but realistic) transformations on our pre-existing training data to create new data. Some of these transformations could be simple replacements of words with their synonyms. Since we had created embedding layers that could define the context of a specific word, we could use that to our advantage, by using cosine similarity to find the similar word for replacement. Alternatively, we could use pre-trained classic word embeddings such as word2vec, GloVe and fasttext to perform similarity search. Another sophisticated data augmentation technique to generate additional, synthetic data, is the back-translation method. In this method, we translate the text data to some foreign language such as French and then translate it back to the English language. This could help to generate textual data with different words while preserving the context of the text data. By generating more data, the network will have a better chance of performing better on the validation and test data.

In conclusion, the construction and implementation of our models may have been successful, but didn't yield the results we wanted to, due to the low recall and precision that each model has. Our future work also would be, to implement the data augmentation technique, as well as use pre-trained word embeddings to try reducing the misclassification error. We could also re-adjust some hyperparameters and reconstruct the architecture of each model to better achieve our main goal. Finally, we hope that the present project of our work will help the humanity and consequently the governments, in order to confront this issue regarding the SDG #3 by taking better actions and implementing more accurate methods.

Appendix

i) Members/Roles

- ΑΓΑΠΙΟΥ ΜΑΡΙΟΣ

1. Web Scrapping Script on PubMed site & Data collection
2. Annotation & Curation
3. Data Parsing & Preprocessing & Model Construction & Model Evaluation
4. Report the Methodologies used to create the models

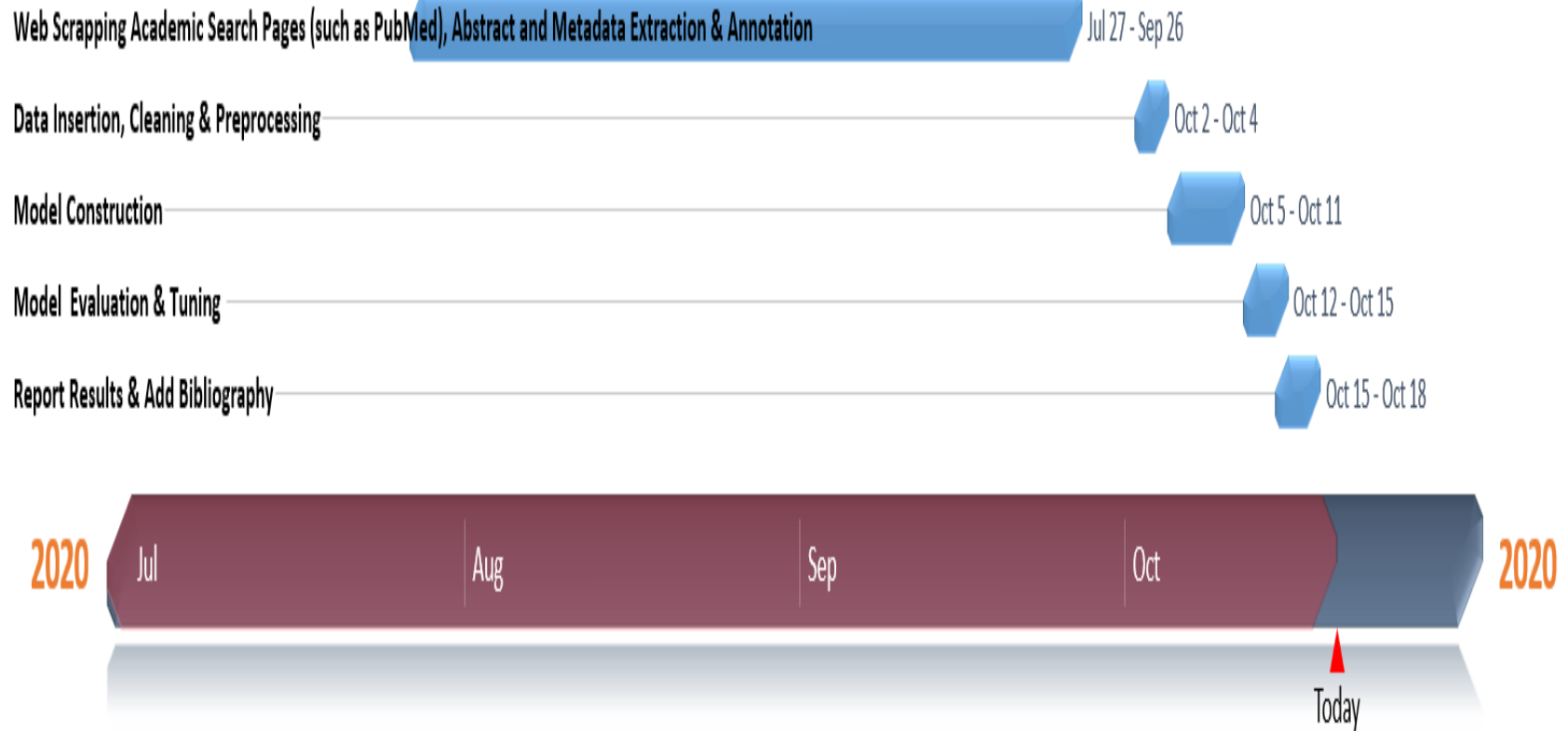
- ΚΟΝΤΟΣ ΧΡΗΣΤΟΣ

1. Data collection
2. Annotation & Curation
3. Data Parsing & Descriptive Statistics & Preprocessing & Model Construction & Model Evaluation
4. Report the Business Problem Description, the Mission and Data Preprocessing

- ΣΚΑΡΛΗΣ ΒΑΣΙΛΗΣ

1. Data collection
2. Annotation
3. Data Parsing & Preprocessing & Model Construction & Model Evaluation
4. Report the Evaluation & Comparison of the models

ii) *Time Plan*



Bibliography

1. Jason Brownlee. A Gentle Introduction to Padding and Stride for Convolutional Neural Networks. machinelearningmastery.com. [Online] April 19, 2019 . <https://machinelearningmastery.com/padding-and-stride-for-convolutional-neural-networks/>.
2. Le Blanc, David, Clovis Freire, and Marjo Vierros. 2017. "Mapping the Linkages between Oceans and Other Sustainable Development Goals: A Preliminary Exploration." DESA Working Paper 149 (February). <https://www.un.org/development/desa/publications/working-paper/wp149>
3. Vladimirova, Katia, and David Le Blanc. 2015. "How Well Are the Links between Education and Other Sustainable development Goals Covered in UN Flagship Reports? A Contribution to the Study of the Science-Policy Interface on Education in the UN system." DESA Working Paper 146 (October). <https://www.un.org/development/desa/publications/working-paper/education-and-sdgs-in-un-flagship-reports>
4. Le Blanc, David. 2015. "Towards Integration at Last? The Sustainable Development goals as a Network of Targets." DESA Working Paper 141 (March). <https://www.un.org/development/desa/publications/working-paper/towards-integration-at-last>
5. Sovrano, Palminari, Vitali. 2004 "Deep Learning Based Multi-Label Text Classification of UNGA Resolutions". <https://arxiv.org/ftp/arxiv/papers/2004/2004.03455.pdf>
6. Rodriguez Medina, Samuel. 2019. "Multi-Label Text Classification with Transfer Learning for Policy Documents: The Case of the Sustainable Development Goals". <http://www.diva-portal.org/smash/record.jsf?pid=diva2%3A1360968&dsid=3657>
7. Chunting Zhou, Chonglin Sun, Zhiyuan Liu, Francis C.M. Lau. 2015. "A C-LSTM Neural Network for Text Classification". https://www.researchgate.net/publication/285458953_A_C-LSTM_Neural_Network_for_Text_Classification
8. Long short-term memory. Wikipedia. [Online] https://en.wikipedia.org/wiki/Long_short-term_memory.
9. ModelCheckpoint. Keras. [Online] https://keras.io/api/callbacks/model_checkpoint/
10. <https://arxiv.org/pdf/1706.00188v1.pdf>
11. <https://paperswithcode.com/paper/a-c-lstm-neural-network-for-text>
12. <https://www.ijitee.org/wp-content/uploads/papers/v8i6s/F60970486S19.pdf>
13. http://www.ijaerd.com/papers/finished_papers/Sentiment%20Analysis%20with%20Multilayer%20Perceptron%20using%20Emoticon%20Space%20Model-IJAERDV04I0760860P.pdf
14. https://www.researchgate.net/publication/336058961_Affective_Computing_and_Deep_Learning_to_Perform_Sentiment_Analysis
15. https://www.researchgate.net/publication/327682642_Investigation_of_Naive_Bayes_combined_with_Multilayer_Perceptron_for_Arabic_sentiment_analysis_and_opinion_mining

16. <https://arxiv.org/abs/1408.5882>
17. <https://paperswithcode.com/paper/convolutional-neural-networks-for-sentence>
18. <https://ai.stanford.edu/~ang/papers/ICPR12-TextRecognitionConvNeuralNets.pdf>
19. <https://papers.nips.cc/paper/5782-character-level-convolutional-networks-for-text-classification.pdf>
20. <https://paperswithcode.com/paper/an-end-to-end-trainable-neural-network-for>
21. <https://paperswithcode.com/paper/multi-branch-attentive-transformer>
22. <https://paperswithcode.com/paper/sentence-state-lstm-for-text-representation>
23. <https://paperswithcode.com/paper/revisiting-lstm-networks-for-semi-supervised-1>
24. <https://arxiv.org/pdf/1901.06610v2.pdf>
25. http://jens-lehmann.org/files/2019/epia_simple_lstm.pdf
26. <https://arxiv.org/pdf/1810.03660v1.pdf>
27. https://en.wikipedia.org/wiki/Feedforward_neural_network
28. https://en.wikipedia.org/wiki/Convolutional_neural_network
29. https://en.wikipedia.org/wiki/Recurrent_neural_network
30. https://www.researchgate.net/publication/43121576_A_Review_of_Machine_Learning_Algorithms_for_Text-Documents_Classification
31. https://www.researchgate.net/publication/228084521_Text_Classification_Using_Machine_Learning_Techniques
32. <https://repositorio-aberto.up.pt/bitstream/10216/107762/2/219435.pdf>
33. <https://medium.com/@juskys8/text-classification-using-machine-learning-5838293c9f2>
34. <https://arxiv.org/pdf/2004.03705.pdf>
35. <https://arxiv.org/pdf/1810.03660v1.pdf>
36. <https://zindi.africa/competitions/sustainable-development-goals-sdgs-text-classification-challenge>
37. https://www.un.org/esa/desa/papers/2019/wp159_2019.pdf
38. <https://www.scitepress.org/Papers/2018/67305/67305.pdf>
39. <https://academic.oup.com/database/article/doi/10.1093/database/baz034/5424138>
40. <https://unctad.org/system/files/non-official-document/PPT%20for%20UNCTAD%20Stats%20-%20by%20Richard%20Rothenberg%20-Feb03-2020.pdf>
41. <https://alexmoltau.medium.com/undp-the-sdgs-and-the-role-of-innovation-involving-machine-learning-ae4ce920d592>
42. <https://www.ijcai.org/Proceedings/15/Papers/033.pdf>
43. <https://inception-project.github.io/>
44. <https://www.clarin.gr/en/about/what-is-clarin>
45. <https://sdgs.un.org/goals>
46. <https://www.un.org/development/desa/publications/working-paper/education-and-sdgs-in-un-flagship-reports>
47. <http://www.diva-portal.org/smash/record.jsf?pid=diva2%3A1360968&dswid=3657>
48. <https://en.wikipedia.org/wiki/AlexNet>
49. https://en.wikipedia.org/wiki/State_of_the_art