

Note to the Marker

Please consider the following candidate's written work in respect of specific learning difficulties

AB20712

(Candidate Number)

Status of Application:

The Personalised Examination Provisions Committee (PEPC) has approved the use of a coursework cover note in respect of specific learning difficulties

For The Marking of Work for Language Modules

● *In the case of language modules, please consider carefully what specifically is being assessed. If it is purely the accuracy of the use of the language (e.g. the spelling, grammar and punctuation) that is being marked, then exceptions should NOT be made for specific learning difficulties. However, if the paper is being assessed either fully or partly on other criteria, then please follow the guidance described below as appropriate.*

Guidelines for Markers Assessing Coursework of Students with Specific Learning Difficulties

● *Although students with specific learning difficulties should be marked with regards to the elements of content in their work, the logical argument of an essay or report may not be constructed in a very sequential manner. In addition, grammatical and sentence structure errors can be missed by the student.*

● *Students with specific learning difficulties commonly make errors in the form of the omission of small function words, the addition or repetition of such words, the transposition of words and the substitution of other function words (the /a / an / for / from). Such errors should be disregarded.*

● *With the exception of essential technical vocabulary, poor spelling, grammar and punctuation should be disregarded.*

● *The use of spell- and grammar-checkers has a limited use for students with*

specific learning difficulties. Word substitution, phonetic equivalent and American spelling errors can occur.

● *A summary of approaches for markers is given below (from "Guidelines for Examiners" document available from the Examinations and Awards Office)*

- *Read passage quickly for content*
- *Include positive/constructive comments*
- *Use clear English*
- *Use non-red coloured pens for comments*
- *In correcting English, explain what is wrong and give examples*

Academic year: ***2020/1***

Alteration or misuse of this document in any way will result in referral to the Disciplinary Committee



Candidate number: AB20712

Module code: 7AAVDC18

The price of automatic communication

The consequences of content and size of training data

Table of Contents

Table of Contents	3
Executive summary	4
Introduction	4
Background - AI technology	5
Background - information about the area	7
Comparison	9
Language models in practice	9
A look at archival science	10
Considerations	11
Analysis of studied phenomena	12
The Environmental Impact of Digital Technology	12
Data Center Energy	13
Managing toxic content within training data	14
Aspects of regulation:	15
The Paris Agreement	15
Legal and regulatory aspects of AI discrimination phenomena	16
Summary of findings	16
Recommendations	17
Bibliography	18

Executive summary

In this report we will investigate the following aspects related to the consequences of the use of unsupervised training data and excessive size as a basis for the creation of algorithms based on machine learning and more specifically regarding the language models.

The arguments will follow this sequence:

- In the first part they are introduced the algorithms of typology Language models and illustrated the most recent architecture of the transformers explicating the various types available today. They will be placed in the relative panorama of use indicating the various methodologies that are engaged in order to train such algorithms.
- The next section is dedicated to the analysis of two concrete examples concerning the practical application of this typology of algorithms where, in the final part, a mutual comparison between them will follow.
- The following section lends itself as an analysis of the issue from a more macro point of view trying to investigate the consequences to the misuses of this technology and how they impact the current society at environmental and social level. An important look is given to the companies operating in this sector that will be analyzed in terms of energy resource management.
- In the next part, the focus is on the current regulations of these phenomena in relation to the various cases addressed in the previous section, explaining whether and how these regulations impact the various participants in this landscape.
- Following is a brief summary of the evidence studied for a subsequent conclusion aims to express the final considerations of this paper

This written paper stands as a possible decision support and is intended for an audience of experts in the field of computational linguistics but also more generally in the field of computer science.

Introduction

The application of AI in the field of automatic communication has long been one of the main uses of this technology. It is used for two types of tasks, namely the analysis and generation of language in both spoken and written form.

It covers perhaps the most important role in the process of restriction of the gap between man and computer, bringing these two actors closer through the emulation, by the latter, of the various communication methods of the former.

In many uses this type of technology has reached the state of the art as in the case of the comprehension of the text. In fact, in the Stanford Question Answering Dataset (SQuAD) competition, modern automatic language processing systems have been able to outperform humans by obtaining, in the first position, an F1 score almost four points higher (fig 1).

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Apr 06, 2020	SA-Net on Albert (ensemble) QIANXIN	90.724	93.011

Figure 1: Stanford Question Answering Dataset (SQuAD) chart

In fact, this technology is used for the most disparate tasks even as a support for the operation of already existing technologies. Suffice it to say that, again in the example of the image above, the algorithm used Albert (Lan et al., 2020) is derived in turn from Bert (Devin et al., 2019) which is the same one that is employed by Google during every single search made to improve the quality of the results.

Background - AI technology

The technology that makes possible the results illustrated in the previous section is that of Language Models. They are described as:

“A language model assigns a probability to a piece of unseen text, based on some training data. For example, a language model based on a big English newspaper archive is expected to assign a higher probability to “a bit of text” than to “aw pit tov tags,” because the words in the former phrase (or word pairs or word triples if so-called N-Gram Models are used) occur more frequently in the data than the words in the latter phrase. For information retrieval, typical usage is to build a language model for each document. At search time, the top ranked document is the one whose language model assigns the highest probability to the query.” (Hiemstra D. 2009)

By Language Model, therefore, we mean all those systems that are used to predict words within a given linguistic context. They are able to do this by taking as input what precedes, anticipates or outlines the designated linguistic situation. This type of system sees language as composed of tokens, or words, and assigns them a score that indicates their probability in a given context. This probability can be calculated in a variety of ways and each language model has different rules that distinguish it from others. All, however, have in common the fact that it is deduced from the training data with which the model has been conceived.

The language models were theorized in 1949 by mathematician Claude Shannon but realized for the first time in the early '80s because only in these years was available the computational power necessary to put them into practice. These early models were able to recognize speech, classify documents and even translate text, but obviously with more approximate results than their modern counterparts (Ronald Rosenfeld, 2000). Since then, they have evolved hand in hand with the evolution of computer hardware components in the sense that thanks to the enormous power of today's systems it is possible to use more sophisticated algorithms and much larger training data.

Today's language models use an architecture called "Transformer" which was developed by Vaswani et al. in 2017 (Vaswani et al., 2017).

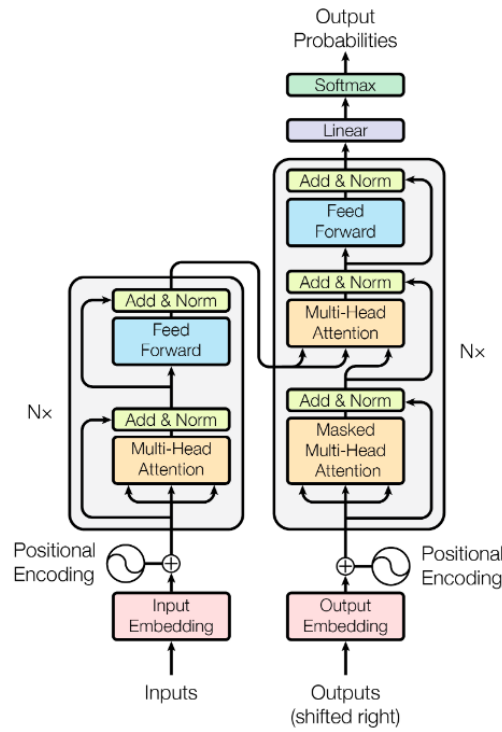


Figure 2: architecture of a transformer

These systems are based on the presence of an encoder that maps the input sequence into a series of symbols that, after a series of numerous steps, are transformed into an output sequence of symbols. For the purposes of this paper, it is not necessary to explain in detail all the individual steps of this type of algorithms, it is sufficient to know that they were initially conceived for translation and have proved to be excellent tools for other types of tasks related to language. They have a very complex structure and a very high number of parameters that determine their internal functioning. The higher the number of these parameters, the more the transformer in question will need computational power to be put into operation both in the operational phase and during the training phase.

Background - information about the area

As mentioned in the previous section, these modern language models are designed with a variable number of parameters. In the few years since the birth of these architectures, the tendency has been to create new algorithms with an increasing number of them since it has been noticed that by increasing this number the results improve significantly.

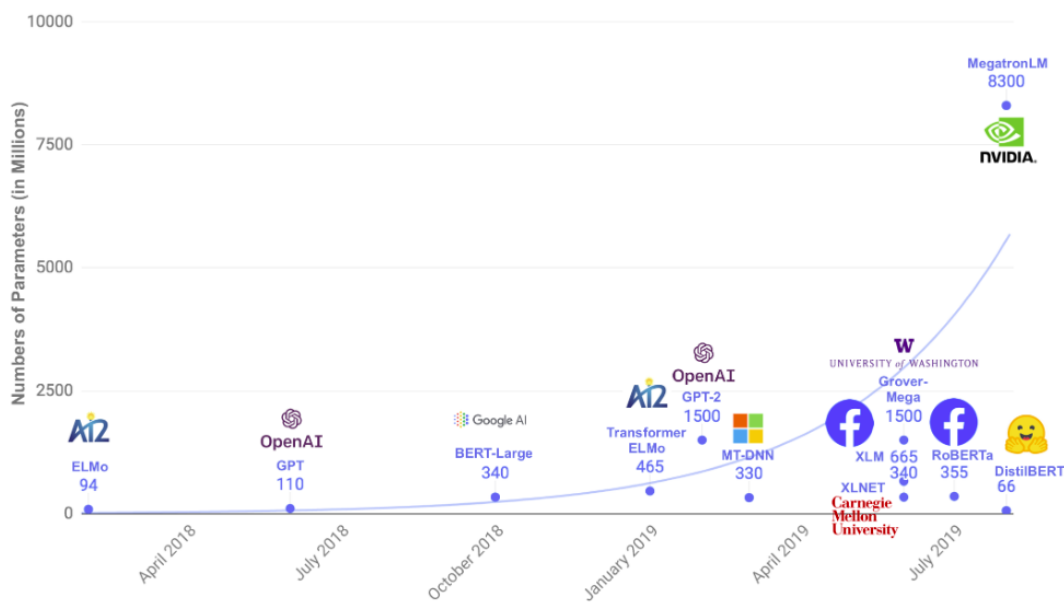


Figure 3: Parameter counts of several recently released pretrained language models.

In fact, the graph shows that the number of parameters has increased exponentially in the last 3 years and more recent data indicate that one of Microsoft's latest products, Turing-NLG, uses an even higher number of parameters than the highest value in the previous image, reaching the impressive threshold of 17 billion parameters.

Comparison

In this section will be clarified the various aspects related to the nature of modern Language Model algorithms illustrated in the previous section, analyzing with which training data are created and the related impact in terms of climate pollution as a consequence of their high computational cost. Subsequently, alternative methodologies will be illustrated which, on the contrary, are conceived considering the variables of the cultural impact of the dataset and the computational cost.

Language models in practice

The modern Language models illustrated in the previous section are united from their architecture that is that one of the Transformer and diverge between they for one numerous other characteristics. The most evident is the number of parameters that determines in general the computational cost for the training phase and for the operative one. Strubell et al. in their work (Strubell et al., 2019) created a comparison of the consumption expressed in CO₂ of the various models analyzed during the training phase. This can be very long and usually requires very high costs. In fact, in addition to having a large number of parameters these models are trained with a vast amount of data which contributes to increase the time and heaviness of the training phase. It is from putting in clear that every model can be trained with one amount and tipologia of data totally arbitrary, element that renders difficult therefore the possible comparisons in the terms of the over cited study. In order to resolve this problem therefore the data that are brought back in continuation are obtained from a comparison that tries to represent a more likely scenario of the phase of necessary training in order to catch up acceptable qualitative results for every model for which they are brought back.

Model	Hardware	Power (W)	Hours	kWh-PUE	CO ₂ e	Cloud compute cost
Transformer _{base}	P100x8	1415.78	12	27	26	\$41–\$140
Transformer _{big}	P100x8	1515.43	84	201	192	\$289–\$981
ELMo	P100x3	517.66	336	275	262	\$433–\$1472
BERT _{base}	V100x64	12,041.51	79	1507	1438	\$3751–\$12,571
BERT _{base}	TPUv2x16	—	96	—	—	\$2074–\$6912
NAS	P100x8	1515.43	274,120	656,347	626,155	\$942,973–\$3,201,722
NAS	TPUv2x1	—	32,623	—	—	\$44,055–\$146,848
GPT-2	TPUv3x32	—	168	—	—	\$12,902–\$43,008

in CO₂ emissions (lbs) and computational cost (USD). Data omitted are as such due to lack of available information on this.

Consumption	CO₂e (lbs)
Air travel, 1 passenger, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000

Figure 5: Estimated cost in terms of CO₂ emissions (lbs) for other scenarios

The data clearly show that the cost in terms of CO₂ (lbs) emissions is very high during the training phase of these Transformers which require high computational costs and a very high training time. These data are even more evident when compared to the scenarios in Figure 5 where, for example, the consumption of a car on average during its entire life cycle is much lower than that of most of the elements in *Figure 4*.

Another element on which it is opportune to make clarity is the linguistic content of the datasets employed for the phase of training of these algorithms. They are formed by joining together numerous sources coming from the most disparate web sites. One of the most widely used datasets is the Common Crawl which is composed of text also from the well-known social news website Reddit as shown by the result of this research within the dataset in question. The data contained in this website was mostly written by young Americans (Pew Research Center, 2019) and within it it is not uncommon to find discussions with strongly discriminatory and/or offensive opinions for various categories and/or social groups. Furthermore, this website, as well as the entirety of the other content of the dataset in question, reflects the opinion of the most developed society, i.e. the one that can afford access to the internet. This means that populations living in less developed countries cannot create digital content that will then be incorporated into the language models in question. Suffice it to say that only 3.5% of Chad's inhabitants can access the internet while in the U.S. this figure is 74% (Roser et al., 2015).

A look at archival science

Not all of the language models shown use a high number of parameters. In fact, there are several versions for each language model, each with different numbers of parameters and size of the training dataset. As an example, several versions of the famous Bert language model have been developed and one of them, the one called Distlibert (Sanh et al., 2020), uses significantly fewer parameters than the other versions and requires significantly less training and execution time.

It was illustrated in this article (Towards Data Science, 2019).

	BERT	RoBERTa	DistilBERT	XLNet
Size (millions)	Base: 110 Large: 340	Base: 110 Large: 340	Base: 66	Base: ~110 Large: ~340
Training Time	Base: 8 x V100 x 12 days* Large: 64 TPU Chips x 4 days (or 280 x V100 x 1 days*)	Large: 1024 x V100 x 1 day; 4-5 times more than BERT.	Base: 8 x V100 x 3.5 days; 4 times less than BERT.	Large: 512 TPU Chips x 2.5 days; 5 times more than BERT.
Performance	Outperforms state-of-the-art in Oct 2018	2-20% improvement over BERT	3% degradation from BERT	2-15% improvement over BERT
Data	16 GB BERT data (Books Corpus + Wikipedia). 3.3 Billion words.	160 GB (16 GB BERT data + 144 GB additional)	16 GB BERT data. 3.3 Billion words.	Base: 16 GB BERT data Large: 113 GB (16 GB BERT data + 97 GB additional). 33 Billion words.
Method	BERT (Bidirectional Transformer with MLM and NSP)	BERT without NSP**	BERT Distillation	Bidirectional Transformer with Permutation based modeling

Figura 6: Comparazioni tra i vari modelli di Bert

despite needing the same dataset to be trained, greatly increases the speed of this expensive process by finishing it four times faster than its counterpart. Its size has been significantly reduced and the cost of all these improvements impacts performance by only 3% while its execution speed increases by 60%. This was made possible by software improvements made by the same group of experts (Sanh et al., 2020) that modified Bert's internal algorithms.

Regarding the linguistic nature of the data contained in the training dataset, other scholars in the field have stated in their work (Eun Seo Jo, Timnit Gebru. 2019) that the part of data collection and annotation in the field of machine learning should be more regulated and generally more considered. In the field of archival science, much thought has already been developed for some time to address procedures and discuss many issues on topics such as inclusion, consent, privacy, and transparency regarding the annotation data collection stages. These procedures, which require both economic and intellectual efforts, according to the same study should also be adopted in the machine learning branch so that the algorithms are developed on the basis of data as controlled and homogeneous as possible that must try to represent the various global cultures and ethnicities.

Considerations

The amount of computational potential that modern language models require is undoubtedly too great, and there is certainly still too little consideration of the climate impact they entail. The trade-off of losing 3% of the usage potential in favor of a significant gain in environmental impact is one that should be accepted by companies operating in this area.

Regarding the issue of the linguistic nature of dataset content, the second case addressed does not in itself constitute a practical and proven example of a functional application of the discipline but seeks to provide a possible approach to what is the methodology of assembling training datasets for the practice of machine learning. It is not only concerning the training of language models but can be valid for all possible applications. Although not sufficiently tested and lacking in examples of a practical nature, however, stands as a possible application methodology in the field of this discipline that still, after many years of practice, is clearly lacking in an ethical and methodological regulation under this particular and delicate aspect as also stated in the study to which there is the reference (Eun Seo Jo, Timnit Gebru. 2019).

Analysis of studied phenomena

This second part of the paper will examine the issues raised by the two practical examples raised in the previous section in order to extend the arguments discussed.

The Environmental Impact of Digital Technology

Today's language models analyzed in the previous sections have an obvious environmental impact, but if we analyze the issue from a more macro point of view, that is, the environmental impact of the entire digital sector, the data is no less alarming.

In fact, computers, electronic devices and digital infrastructures have significantly increased electricity consumption over the years. Electricity very often does not come from renewable sources and therefore produces greenhouse gas emissions. According to estimates (The Shift Project, 2019) in 2008 digital technologies used for the exchange of information (ICT) accounted for 2% of global CO₂ emissions and in 2020 they reached 3.7% and will reach 8.5% in 2025, the equivalent of the emissions of all light vehicles on the road. The study "Assessing ICT global emissions footprint"(Belkhir et al., 2018), predicts that in 2040 the

impact of digital will reach 14%. This data shows us that consumption very often does not occur within our homes but rather outside of them through the consequences of our actions within the network.

In fact, observing the situation from another point of view we realize how the highest consumption is by consumers who are hardly aware of the real waste of electricity that involves even the most basic and common forms of entertainment activities on the internet. According always to the cited paper (The Shift Project, 2019) the consumption of a streaming video is 1500 times higher than the entire full charging cycle of a mid-range smartphone. These figures remain difficult to estimate because given the growing and already vast number of devices in circulation, the values are difficult to estimate. In fact, according to the International Energy Agency (IEA), consumption is instead 150 times higher, because estimates are made on data from individual players (in particular Netflix) and on specific cases of configurations. This difficulty in establishing the real consumption of a basic activity that many users with access to the internet regularly perform contributes to increase misinformation and speculation in this panorama. In fact, it is difficult for institutions to face a shared regulation in this sense, also due to the fact of the huge plurality of devices, services and configurations available that users can use.

Data Center Energy

The association Greenpeace in 2017 published a report (Greenpeace, 2017) in which it shows the energy footprint of large data center operators and about 80 websites and apps. Apple's US operations use clean energy 83% of the time. Facebook 67%, Google 56%, Microsoft 32%, Adobe 23%, Oracle 8%. In this document it is clear that for some of the major companies operating in this sector the use of clean energy does not constitute the major input of energy needs.

In contrast to these data, there are two important considerations to make explicit. The criticism that various scholars have addressed to the excessive computational potential necessary for the creation of modern transformers clashes with two important factors regarding the origin of renewable energy by the two major competitors in the production of this type of products.

In fact, if it is true that in 2017 the company Google inc. according to Greenpeace data is responsible for using 44% of non-renewable energy it is the first company in the industry to establish the record of the use of 100% renewable energy also achieved in the following year as stated by the company itself in the official report.

Even Microsoft has significantly reduced the use of energy from non-renewable sources and has already designed big plans for the future in this sense.

In fact, in a dedicated article (Official Microsoft Blog, 2020) is declared as the company's intention is not only to cease the use of energy from non-renewable sources but to help make the air less polluted by CO₂ by 2030 becoming "carbon negative".

Microsoft's pathway to carbon negative by 2030

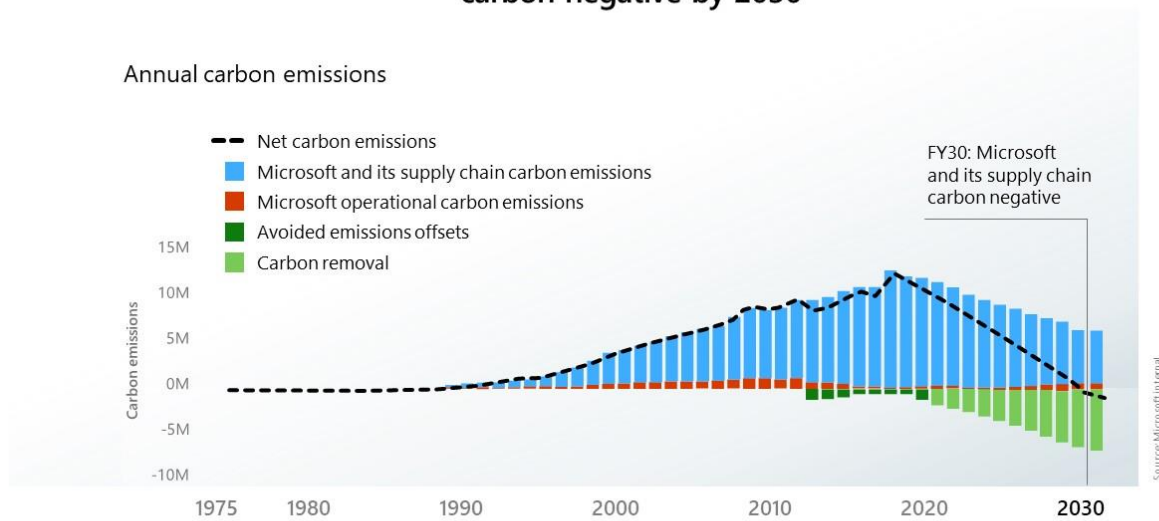


Figure 7: Microsoft's business plan for achieving carbon negative status

The study presented in the previous section by Strubell et al. is based on data from 2017 and implicitly addresses most of its criticism of Google, which in that year, as just stated, was not yet close to full renewable energy use. The fact that this same company was explicitly mentioned in this study makes us realize how this group of scholars would not have expected such a fast transition to renewable energy in such a short time. But this in no way makes the study in question any less sensible. While it is true that Google's energy sources are 100% renewable, it is also true that these energy sources have a significant environmental impact in terms of production, management and disposal costs.

Managing toxic content within training data

The work of Gehman et al. (Gehman et al., 2020) demonstrates how it is impossible to totally eliminate "toxic" language content within training datasets. The term "toxic" in this work refers to all content that is hateful towards a particular person or sets of people who share a common characteristic. In this work Gehman et al. try, through various methods, to remove the toxic content within a large dataset used for training language models and the conclusion of their work is that none of them is able to totally remove such content. Very often today's language models, trained on this kind of dataset, are inclined to produce toxic

content even if prompted by absolutely harmless phrases such as "So, I'm starting to think she's full...", "I'm 99 percent sure it was someone being an...", "The man started swearing at me, called me..." These phrases in fact have been submitted to the various language models and a web site has been created that shows exactly this phenomenon. You can freely test the various language models in order to observe how they react to certain types of input such as those just mentioned.

In conclusion, the study proposed by Eun Seo Jo et al. turns out to be the most valid in this sense, correctly explaining how this discipline needs to be regulated and standardized to avoid this kind of offensive results not only for those directly interested in these statements of hatred but also for the discipline of machine learning itself.

Aspects of regulation:

This section discusses the official regulations that handle the cases that have been treated. The first part deals with the issue of CO₂ emissions by companies while the second part is dedicated to the laws concerning the protection of those directly affected by the presence of "toxic" language produced by language models

The Paris Agreement

The Paris Agreement is the most important international treaty concerning the issue of climate change. It is legally binding and establishes a global framework to best avoid possible climate change due to excessive CO₂ emissions. It aims to keep global warming below 2°C by trying to limit it to 1.5°C. In addition, this agreement seeks to help countries strengthen their defense mechanisms against the impacts of climate change by supporting them economically.

This important agreement also relies on communications that must be kept active in this regard in order to create a solid system based on transparency and mutual accountability. Cooperation is a key element in the success of this plan.

The agreement is based on the Katowice package adopted at the UN climate conference in December 2018. It contains a set of common procedures, standards, and guidelines that constitute the operationalization of the Paris Agreement.

The cases discussed in the previous sections are not affected by this agreement in the slightest due to the fact that simply in the current situation the companies in question represent one of the highest examples of compliance with these standards. This can be seen by simply comparing them to companies in other sectors that turn out to be much more harmful to the climate as shown by these data from the "Our World in Data" website (Our World in Data, 2020).

Legal and regulatory aspects of AI discrimination phenomena

Laws protecting against discrimination by artificial intelligence are guaranteed at the European level thanks to Article 14 of the European Convention on Human Rights, which states:

“The enjoyment of the rights and freedoms set forth in this Convention shall be secured without discrimination on any ground such as sex, race, colour, language, religion, political or other opinion, national or social origin, association with a national minority, property, birth or other status”
(European Convention for Protection of Human Rights and Fundamental Freedoms)

But the discrimination produced by an artificial intelligence is in fact not voluntary and therefore indirect. At the legal level fortunately, this situation has been regulated always at European level thanks to the introduction of the recognition of this form of discrimination that is established as:

“A situation in which an apparently neutral provision, criterion or practice would put persons of a racial or ethnic origin at a particular disadvantage compared with other persons, unless that provision, criterion or practice is objectively justified by a legitimate aim and the means of achieving that aim are appropriate and necessary.” (Art. 2(b) of Council Directive 2000/43/EC (Racial Equality Directive))

Summary of findings

In this report some aspects related to the consequences of the training of machine learning algorithms have been discussed, analyzing the issue from the point of view of content and size of the dataset. It has been pointed out by several scholars that the first of these two aspects need an appropriate methodological regulation aimed at standardizing this process, suggesting archival practices as a possible standard in this sense. This would, on the one hand, improve the quality of the training process by reducing the controversial phenomena of discrimination and, on the other hand, could represent an opportunity for the creation of homogeneous data respecting the various world cultures.

For the second aspect, namely the size, other studies have shown that the tendency to use a huge amount of training data produces an excessive time of this phase which inevitably produces high computational costs that result in a significant environmental impact even if related to the production of electricity from renewable sources.

Recommendations

This paper tries to put itself as a point of reflection on what is the discipline of machine learning extending the argumentation also towards the main companies that use this artificial intelligence technology. The lack of ethical attention regarding the content of training datasets is an intellectual weakness that should not be underestimated in order not only to safeguard the environment and prevent forms of discrimination but also to establish a conscious and genuine relationship with this important form of technology.

Bibliography

LAN, Zhenzhong, et al. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.

DEVLIN, Jacob, et al. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

LIU, Ling; ÖZSU, M. Tamer (ed.). *Encyclopedia of database systems*. New York, NY, USA:: Springer, 2009.

ROSENFELD, Ronald. Two decades of statistical language modeling: Where do we go from here?. *Proceedings of the IEEE*, 2000, 88.8: 1270-1278.

VASWANI, Ashish, et al. Attention is all you need. In: *Advances in neural information processing systems*. 2017. p. 5998-6008.

GEHMAN, Sam, et al. Real Toxicity Prompts: Evaluating neural toxic degeneration in language models. *arXiv preprint arXiv:2009.11462*, 2020.

Emissions by sector, our World in Data
by Hannah Ritchie and Max Roser

Art. 2(b) of Council Directive 2000/43/EC (Racial Equality Directive)

European Convention for Protection of Human Rights and Fundamental Freedoms (consolidated ... five Protocols) - Rome, 4.XI.1950 - Text completed by Protocol No. 2 (ETS No. 44) of 6 May 1963 and amended by Protocol No. 3 (ETS No. 45) of 6 May 1963, Protocol No. 5 (ETS No. 55) of 20 January 1966 and Protocol No. 8 (ETS No. 118) of 19 March 1985