## Statistics on companies undertaking activities in the field of corporate social responsibility (CSR) using web scraping and machine learning

Organisation: Statistics Poland

Authors: Bartosz Grancow,
Emilia Murawska,
Klaudia Peszat

Date: 15.11.2022

Version: 2.0

### Introduction

The traditional data sources do not provide information on many new and fast changing socio-economic phenomena, therefore national statistical offices more often reach for new data sources, such as web data.

Statistics Poland within the ML2022 web scraping theme group made an attempt to explore web data sources and new machine learning techniques to recognize the possibilities to augment official statistics on enterprises. The experimental research was focused on Corporate Social Responsibility activities.

Corporate Social Responsibility (CSR) is defined as „a company's sense of responsibility towards the community and environment (both ecological and social) in which it operates.

Companies express this citizenship:
    (1) through their waste and pollution reduction processes,
    (2) by contributing educational and social programs, and
    (3) by earning adequate returns on the employed resources" (*the Business Dictionary).*

The report presents the results of the research which aimed at the classification of companies into two groups: those which carry out activities in a CSR field (and inform about it) and those which more likely do not have any CSR strategies or at least do not communicate their activities in this area. The type of such activity was not the subject of the project.

### Data

#### Input data

The data source used to build the population of companies was the Business Register – Database of Statistical Units (BJS), including name of enterprise, address, e-mail address, URL address (if available). Due to the fact that CSR-related activities are most often carried out by large and medium-sized companies (small companies often do not have websites), we focused our attention on this group of enterprises (employing more than 10 employees).

The first step was the extraction from the Business Register of the population on enterprises with URL addresses. It is worth mentioning that, the URLs can be also obtained from email address. The methodology for obtaining URLs among others from email addresses has been described in detail in the report delivered within the ESSnet Trusted Smart Statistics – Web Intelligence Network project (Kuhnemann et al. 2022). Taking into account the exploratory purpose of our research and the fact that the database of large and medium companies covered over 20 thous. enterprises with URL addresses, we did not decide to obtain additional URLs via email addresses. The total number of URLs in our database consisted of 6 thous. URLs for large enterprises and nearly 14 thous. for medium ones from all over Poland.

Next, we checked validity of the URLs using web scraping tools. There were 11,380 valid URLs (with response code 200). Inactive URLs were removed from the database.

**Search engine results**

We have decided that in a scope of our interests are only websites held by the companies to provide information on their business. The enterprise should have control over the content of such websites. External web pages providing basic information on the companies, such as yellow pages, government domains etc. were excluded.

Information on CSR activities can be found in various places on the website. Some enterprises post them on the home page, others have a dedicated tab or subpage. The descriptions of the CSR activities conducted by companies are also published in various ways. Some enterprises publish them as a plain text on the website, others attach pdf reports.

Taking into consideration the enormous variety of websites, we decided to use an internet search engine to find texts from the official company's web sites referring to their activities in the CSR field. Only search results (snippets) were used as data source and no additional scraping was done. This approach has been implemented first by CBS (Delden et al. 2019). This enabled to save time and resources needed to web scrape the entire content of webpages.

The choice of the search engine and the search term have significant impact on the results. However, due to the exploratory purpose of the research and payment requirements for using popular search engines, such as Google or Bing for mass queries, we decided to use the less known search engine – Duck Duck Go[1]. Its advantage is not collecting data about users, and thus not profiling the results. Each user can get the same search results, which made it easier to work on the project.

We also tested various search options by entering company's URL and keyword or name of a company and keyword. The first variant brought more accurate results. When entering a company's name, contact details were usually in the first places of the search.

To prepare a catalogue of keywords that appeared most often in the descriptions of CSR activities we also searched manually on companies websites. The directory consisted of words and phrases such as: 'csr', 'społeczna odpowiedzialność', 'odpowiedzialność biznesu', 'społeczna', 'zaangażowanie', 'środowisko', 'ekologiczne', 'zrównoważony rozwój', 'csr raport', 'charytatywny'. We tested all above mentioned keywords, however the search term 'URL address' + keyword: 'csr' gave the best results.

We also noticed that if a company conducts CSR activities, information about it appears in the first five search results (most often in the first three descriptions, but there are some cases when it appears in the fourth or fifth ones). For this reason, we limited the search to the first five results.

Below is an example of the search results – URL address and keyword entered manually in the search engine Duck Duck Go and the results of descriptions retrieved and saved in an Excel file.
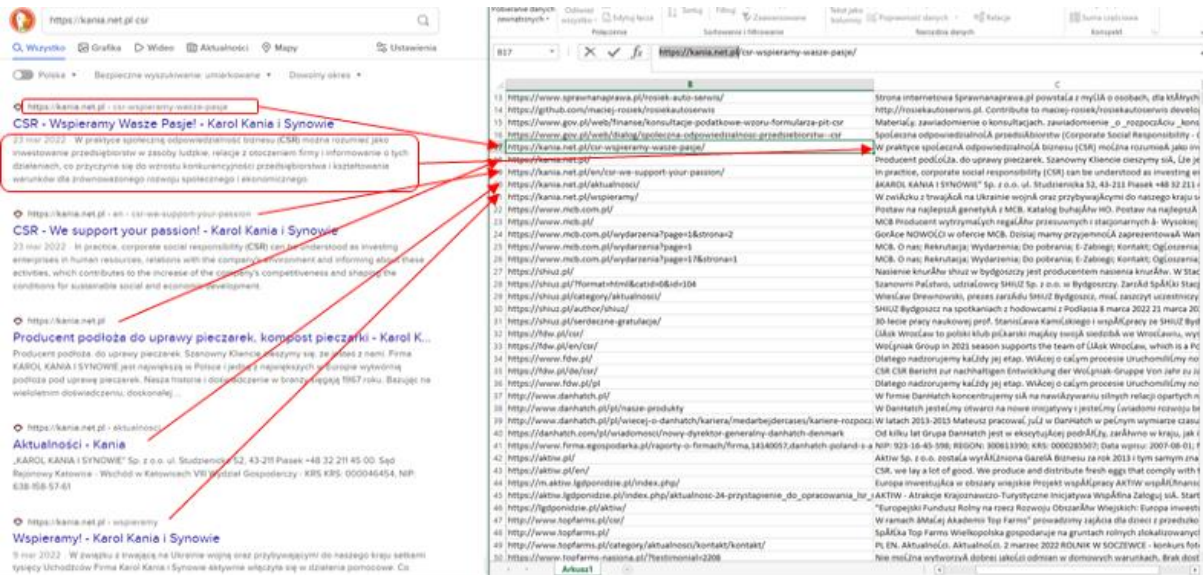
**Figure 1. Search engine results**


## Data preparation

Data stored in an Excel file had to be further processed. At this step, results which did not come from company's' webpages were removed from the database. It was: yellow pages, phonebook websites (e.g. panoramafirm.pl, oferteo.pl, fimy-pl.com), job offer pages (e.g. gowork.pl, pulshr.pl), social media (e.g. Facebook, Linkedin), social campaign pages (e.g. kampaniespoleczne.pl), non-governmental organizations pages (addresses with the org.pl domain), as well as government and educational institution pages (addresses with the gov.pl and edu.pl domain).

In the next step, the pre-processing techniques, such as tokenization, reduction of words to their basic form by lemmatization and stop words removal were applied.


## Machine learning solution

### Model Top2Vec

Once the data was prepared, we attempted to detect proximity of texts from our database with the CSR topic. For this purpose, we applied a new unsupervised machine learning algorithm – the Top2Vec. The Top2Vec algorithm is used to automatically detect topics present in texts and documents. This algorithm was chosen as it enables working on short text (such as snippets), does not require removing stop words and pre-processing (e.g. stemming/lemmatization)[2].

Despite the fact that the Top2Vec algorithm does not require previous text processing, some of the pre-processing steps were done, mostly to compare the model results applied on the raw data and pre-processed ones.

The default embedding model applied in the Top2Vec algorithm which enables creation jointly embedded document and word vectors is Doc2vec. This is not a pre-trained model, it learns from scratch and thus it can be used for different languages. The algorithm provides also a possibility to use pre-trained embedding models, such as: universal-sentence-encoder, which currently supports 16 languages, including Polish. For the experimental research purposes different embedding models were tested, however universal-sentence-encoder-multilingual-large model proved to give the best results.

## Software used

The implementation of the Top2Vec model was carried out in a Python programming language. The installation of the library seemed not to be very challenging, however several adjustments of the IT environment were necessary. First of all, Microsoft Visual C++ 14.0 or greater was required, and secondly – installation of hdbscan. The latter problem can be easily solved in Anaconda environment by running the code: conda install -c conda-forge hdbscan.

## Results

### Raw data

The Top2Vec algorithm was applied on the dataset consisting of 13.4 thous. raw texts. The model detected 73 topics, including one very closely related to the CSR concept (Topic 0). It consists of 477 semantically similar texts, which also makes it the largest cluster. The wordcloud for the topic 0 is presented below.
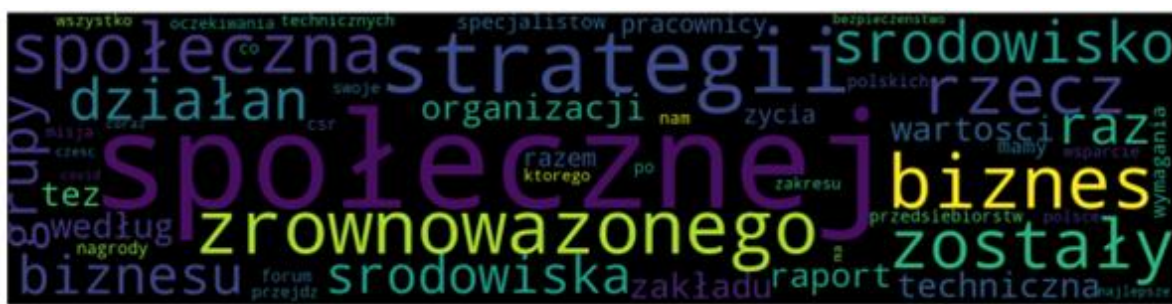


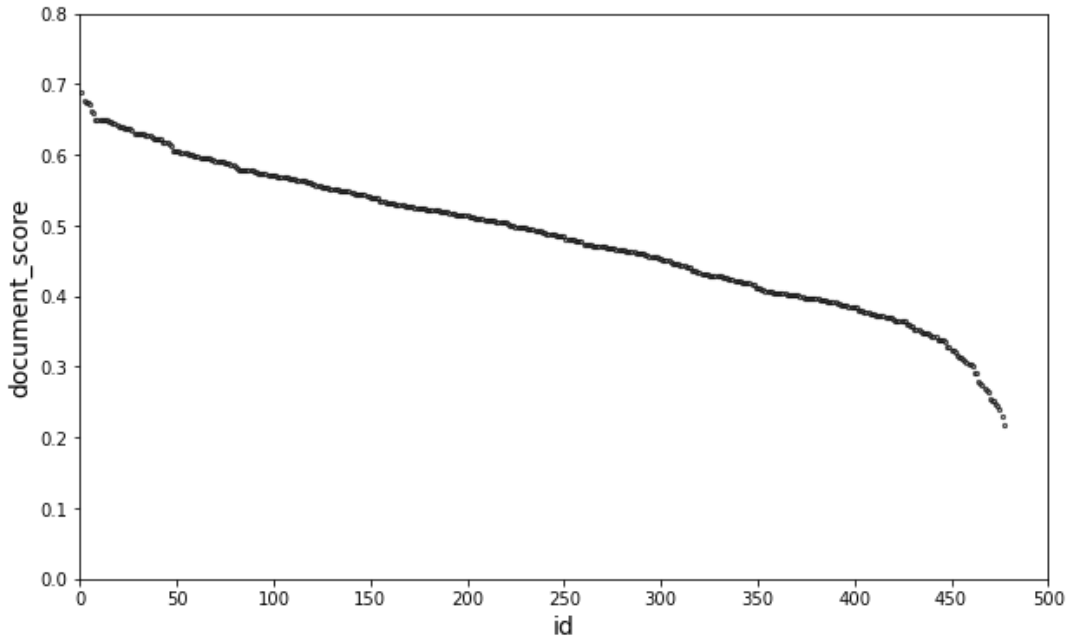**Figure 2. The wordcloud for the topic 0 (in Polish)**

In the next step, we calculated the document scores, which indicate a semantic similarity of documents to the assigned topic. The table below presents the document scores for the texts clustered in the topic 0 (in descending order).

| | id | url | text | topics | topic_scores | document_scores |
|---|---|---|---|---|---|---|
| 0 | 1 | https://kania.net.pl/csr-wspieramy-wasze-pasje/ | w praktyce społeczną odpowiedzialność biznesu ... | 0.0 | 0.519214 | 0.687659 |
| 1 | 2 | https://kania.net.pl/en/csr-we-support-your-pa... | in practice corporate social responsibility cs... | 0.0 | 0.506223 | 0.675753 |
| 2 | 3 | https://kania.net.pl/de/csr-wir-unterstuetzen-... | in der praxis kann die soziale verantwortung v... | 0.0 | 0.521032 | 0.673456 |
| 3 | 4 | https://www.czatkowice.pl/zrownowazony-rozwoj/csr | csr zrównoważony rozwój kopalnia wapienia czat... | 0.0 | 0.470100 | 0.673364 |
| 4 | 5 | https://www.czatkowice.pl/zrownowazony-rozwoj | kwestię csr traktujemy bardzo kompleksowo dlat... | 0.0 | 0.477510 | 0.671222 |
| ... | ... | ... | ... | ... | ... | ... |
| 472 | 473 | https://about.puma.com/en/sustainability | social compliance against modern slavery and h... | 0.0 | 0.436546 | 0.247481 |
| 473 | 474 | https://annual-report.puma.com/2020/en/sustain... | t puma for sustainability targets performance ... | 0.0 | 0.250929 | 0.245022 |
| 474 | 475 | https://www.csrhub.com/CSR_and_sustainability_... | cnh industrial n v description open who uses c... | 0.0 | 0.395547 | 0.239758 |
| 475 | 476 | https://www.dailycsr.com/tags/CNH%20Industrial/ | daily csr companies environment economics poli... | 0.0 | 0.538010 | 0.230428 |
| 476 | 477 | http://joyevent.pl/party/csr-w-naszym-domu | csr corporate social responsibility jest konce... | 0.0 | 0.624386 | 0.217756 |

477 rows × 6 columns

**Figure 3. The results for the topic 0**

The results can also be easily assessed in a scatter plot where the x axis is the document id and y axis is the document score.

**Figure 4. The scatter plot for the topic 0**

The table and scatter plot contain the number of texts from companies websites, thus to calculate the number of companies which more likely conduct CSR activities we had to deduplicate information.
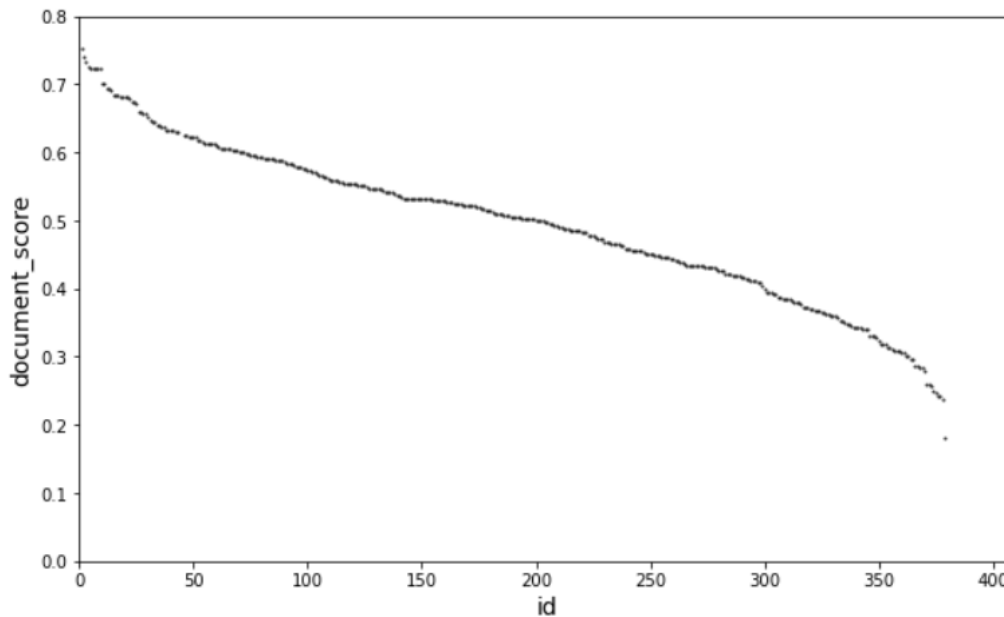
The number of unique companies which are assigned to the topic 0, closely related to the CSR concept is 280 (2,5% of all companies).

**Pre-processed data**

The same process was carried out on the pre-processed data. In this case, the model detected 93 topics, from which topic 2 seemed to be the most related to the CSR concept. It was the third topic taking into account its size (379 documents).



**Figure 5. The wordcloud for the topic 2 (in Polish)**

**Figure 6. The scatter plot for the topic 2**

The number of unique companies which are assigned to the topic 2, closely related to the CSR concept is 246 (2,2% of all companies).

The above means, that the results of the model trained on the raw data and on the pre-processed dataset are very similar. The random assessment of the assignments of texts to the CSR-related topics confirmed these outcomes, however further research as well as more advanced validation methods should be used when considering production of official statistics.

## Reproducibility of the results

Basically, the Top2Vec model is stochastic, thus to achieve the reproducibility of the results for the same dataset it is necessary to set a random seed in the UMAP algorithm. Unfortunately, the Top2Vec library does not provide the possibility to set this parameter via API and it can be only done in the source code[3].

## References

Delden, Arnout van; Windmeijer, Dick; ten Bosch, Olav (2019): *Finding enterprise websites*. Bilbao (European Establishment Statistics Workshop).

Kuhneman, Heidi; van Delden, Arnout; Summa, Donato; Gussenbauer, Johannes; Ils, Alexandra, Loytynoja, Katja (2022): *URL finding methodology*. Joint report for Work Package 2 (Online Based Enterprise Characteristics) and Work Package 3, Use Case 5 (Business register quality enhancement) V.5.0. ESSnet Trusted Smart Statistics – Web Intelligence Network. https://ec.europa.eu/eurostat/cros/content/url-finding-methodology_en