

On the Validity of Using Webpage Texts to Identify the Target Population of a Survey

Citation for published version (APA):

Daas, P., Hassink, W., & Klijs, B. (2024). On the Validity of Using Webpage Texts to Identify the Target Population of a Survey: An Application to Detect Online Platforms. *Journal of Official Statistics*, 40(1), 190-211. <https://doi.org/10.1177/0282423X241235265>

Document license:

CC BY-NC

DOI:

[10.1177/0282423X241235265](https://doi.org/10.1177/0282423X241235265)

Document status and date:

Published: 01/03/2024

Document Version:

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.tue.nl/taverne

Take down policy

If you believe that this document breaches copyright please contact us at:

openaccess@tue.nl

providing details and we will investigate your claim.

On the Validity of Using Webpage Texts to Identify the Target Population of a Survey: An Application to Detect Online Platforms

Journal of Official Statistics
2024, Vol. 40(1) 190–211
© The Author(s) 2024
Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/0282423X241235265
journals.sagepub.com/home/jof



Piet Daas¹, Wolter Hassink^{2,3}, and Bart Klijs⁴

Abstract

A statistical classification model was developed to identify online platform organizations based on the texts on their website. The model was subsequently used to identify all (potential) platform organizations with a website included in the Dutch Business Register. The empirical outcomes of the statistical model were plausible in terms of the words and the bimodal distribution of fitted probabilities, but the results indicated an overestimation of the number of platform organizations. Next, the external validity of the outcomes was investigated through a survey of the organizations that were identified as a platform organization by the statistical classification model. The response by the organizations to the survey confirmed a substantial number of type-I errors. Furthermore, it revealed a positive association between the fitted probability of the text-based classification model and the organization's response to the survey question on being an online platform organization. The survey results indicated that the text-based classification model can be used to obtain a subpopulation of potential platform organizations from the entire population of businesses with a website. This subpopulation may form a good starting point to study platform organizations in more detail.

Keywords

external validation, type-I error, machine learning, web pages

¹Mathematics and Computer Science Department, Eindhoven University of Technology, Eindhoven, The Netherlands

²Utrecht University School of Economics, Utrecht University, Utrecht, The Netherlands

³IZA, Institute of Labour Economics, Bonn, Germany

⁴Sector Culture, Tourism and Technology, Statistics Netherlands, The Hague, The Netherlands

Corresponding author:

Piet Daas, Mathematics and Computer Science Department, Eindhoven University of Technology, De Groene Loper 5, Eindhoven 5612 AZ, The Netherlands.

Email: p.j.h.daas@tue.nl



I. Introduction

Obtaining reliable information from a small or rare subpopulation is a challenging topic for survey researchers (Snijkers et al. 2023; Tourangeau et al. 2014), especially in an era where response rates continue to decline (Luiten et al. 2022; Wu et al. 2022). Approaches commonly used to find rare or so-called hard-to-identify groups are a screening survey, network sampling, area sampling, or a combination (Snijkers et al. 2013). Sometimes, lists of particular types of units are obtained from commercial organizations or they are constructed from administrative data sources (United Nations 2020). Unfortunately, these approaches do not always provide a good overview of the population of interest, especially when the topic of the study is new (Tourangeau et al. 2014; United Nations 2020, chap. 8). However, the increasing availability of new data sources, so-called Big Data (Daas et al. 2015), may offer a solution to this problem. For example, such sources could be used to identify the relevant subpopulation, that is, the target population of the survey, as completely as possible without contacting those units. More specifically, it can be applied to identify businesses with an online platform (defined below).

The surge of internet technology in recent decades has enabled the rapid development of online platform organizations, and it has strongly altered the functioning of society. As a virtual digital meeting point, the intermediary platforms bring together persons and organizations, via which goods, services, or information can be exchanged. The Organisation for Economic Co-operation and Development (OECD 2019) defines online (digital) platforms as “a digital service that facilitates interactions between two or more-distinct but interdependent sets of users (whether firms or individuals) who interact-through the service via the Internet.” Digital labor platforms can, for instance, be applied to a geographically dispersed crowd, for example, “crowdwork,” and by apps (Berg et al. 2018; Howcroft and Bergvall-Kåreborn 2019). Furthermore, online platforms have been applied to shopping (Ducci 2020) as well as to the sharing economy (Sutherland and Jarrahi 2018). A substantial number of online platform organizations, such as Airbnb, Greenwheels, and Uber, are profit-driven, and it may have implications for competition in two-sided markets (Cui et al. 2020; Ducci 2020; Rochet and Tirole 2003).

In recent years, national statistical institutes (NSIs) were lagging behind the phenomenon of the emergence and rapid growth of online platforms. Reliable statistics on the key components and dimensions of the online platform economy are still lacking (United Nations 2019). To fill this gap, in the past years NSIs have debated a framework for measuring elements of the online economy for the Gross Domestic Product and the national accounts (OECD 2020). However, many empirical studies that investigate the (size of the) platform economy are solely based on surveys asking businesses or individuals about their use of online platforms (De Groen et al. 2017). Directly targeting online platform organizations instead of their users has proven to be more difficult (Heerschap et al. 2021; Klijs 2021). The main reason is that the identification of an online platform organization is far from straightforward. Online platform organizations cannot adequately be identified from Business Registers of NSIs, since the business classification system (NACE) classifies businesses according to their main economic activity; the system has no separate category for online platforms (Eurostat 2008). This means that these types of businesses cannot easily be approached with questionnaires assessing

their economic activity. Also, alternative approaches, such as a generic list or a register of online platform organizations, are currently not available.

The inability to obtain a population of online platforms has hindered the direct research of those organizations using questionnaires or administrative data sources. However, because all online platform businesses have a website, it is of interest to investigate if the texts on those websites could be used to accurately identify them. In that way, a list of potential online platform organizations active in a country could be obtained. Text mining techniques could be used to do this (see Becue et al. 2004). In this article, we describe the results of our study which aims to identify the online platform organizations in the Netherlands. Our empirical analysis is based on information obtained from the websites of about six hundred thousand Dutch organizations; these are, in principle, all websites that have been assigned to the businesses in the Business Register of Statistics Netherlands (Daas and Van der Doef 2020). For all these businesses, we have collected the textual content of the pages on their website. During the text mining analysis, we focus on combinations of words that tend to be associated with online platforms. The organizations will be ranked with respect to the likelihood of being an online platform. We demonstrate that the selection of platform organizations obtained is confirmed by an ex-post statistical analysis. The findings are subsequently validated using the information from the Dutch Online Platform survey conducted among the organizations identified as online platforms. To the best of our knowledge, only a limited number of empirical studies have assessed the validity of web-based text mining results through ex-post survey information (García Lozano et al. 2020).

Our study has two major implications. First, it demonstrates that text-based classification is a valid way to obtain a subpopulation strongly enriched with the target population of interest. Second, it demonstrates the advantages of combining text mining techniques and survey data for the study of the online economy.

The setup of this article is as follows. Section 2 describes the general methods used and introduces text mining as a classification method. In Section 3, this method is applied to the texts extracted from the websites of Dutch organizations and the external validity of the text-based classification results is examined using survey information. Finally, in Section 4 the findings are discussed.

2. Material and Methods

2.1. Data Collection and Text Processing

All scripts used are written in Python (v3.7). The Business Register of Statistics Netherlands (Ritzen 2007) is used to provide an overview of all businesses in the Netherlands. To this register, at the most detailed level possible, the corresponding websites are linked. The linking procedure, amongst others, compares the Chamber of Commerce number and address displayed on the website with those in the Business Register; for more details see Oostrom et al. (2016) and Daas and Van der Doef (2020). Websites are assigned to a total of 960,588 organizations, at the level of the local unit. The relationship between these units and their website is essentially one-to-one.

Web pages are collected, that is, scraped, with the `urllib.request` function in Python. Each page is verified with the Beautiful Soup library (v4.7.1) after which it is stored on

the local machine. Pages that can't be scraped during the first attempt are visited at least four times, at later points in time, to deal with temporarily unavailable websites. Scraping started at the main page of the website, followed by all pages referred to on the same website, up to a maximum of one thousand. Collecting all data took 3.5 weeks and resulted in a total of nearly 1 Terabyte of data. The locally stored files are processed in several steps illustrated in Figure 1. The language of the extracted text is determined with the `langdetect` (v1.0.7) library. Since the majority of the pages are either written in Dutch or English only those languages are discerned; for example, any non-Dutch text is classified as English. For the removal of language-dependent stop words the Dutch and English stop word lists in the NLTK-library (v 3.4.1) are used. The words remaining could be stemmed, for example, reducing words to their root form. For this, the `SnowballStemmer` library (v1.2.1) is used. Stemming has the advantage that it considerably reduces the number of variants of a word; for example, the words “helpful,” “helpfully,” and “helping” are all converted to “help”. For model development, either only the text on the main webpage or the texts extracted from all pages collected on the website are used. The texts were combined into a single document in which the words are separated by single spaces. Websites for which ten or fewer words remain after processing, which is particularly relevant when only the text of the main page was studied, are excluded for further analysis as this has been demonstrated to hardly provide any relevant information (Daas and Van der Doef 2020).

To enable model development, the well-known representation of the text extracted in the form of frequency-annotated bag-of-words is used (Aggarwal 2016). This starts by creating a document-term matrix in which the rows correspond to the business webpages and the columns to the unique words included in all the text extracted. The natural logarithm of the term frequency-inverse document frequency ($\log(tf-idf) + 1$) for each word is used as a feature value (Daas and Van der Doef 2020). The *tf-idf* value indicates how important a word is in the texts as the term frequency increases proportionally to the number of times a word appears in it. The inverse document frequency offsets this number by the number of texts that contain the word. The latter adjusts for the fact that some words appear more frequently than others in website texts, of both platform and non-platform texts, which severely reduces the influence of often occurring, non-discriminating, words. In addition, the language of the text is added as a binary feature to the matrix, for which *English*=1 and *Dutch*=0. Word Embeddings, a technique focused on word co-occurrences that is often used to improve text classifications by encoding semantic and syntactic information (Allen and Hospedales 2019), are included by applying the `gensim` library (v3.4.0). Up to three hundred vectors of either the `word2vec` skip-gram or Continuous Bag Of Words algorithms of the `gensim` library could be additionally added to the matrix. Machine Learning models are developed with the `scikit-learn` library (v0.21.2; Pedregosa et al. 2011).

2.2. Model Development and Classification

The overall process of text processing, which is applied in the empirical analyses, consists of three steps. In the first step, a data set with known examples of platform and non-platform websites is constructed by experts. All three experts involved are employees of Statistics Netherlands with at least five years of experience in business

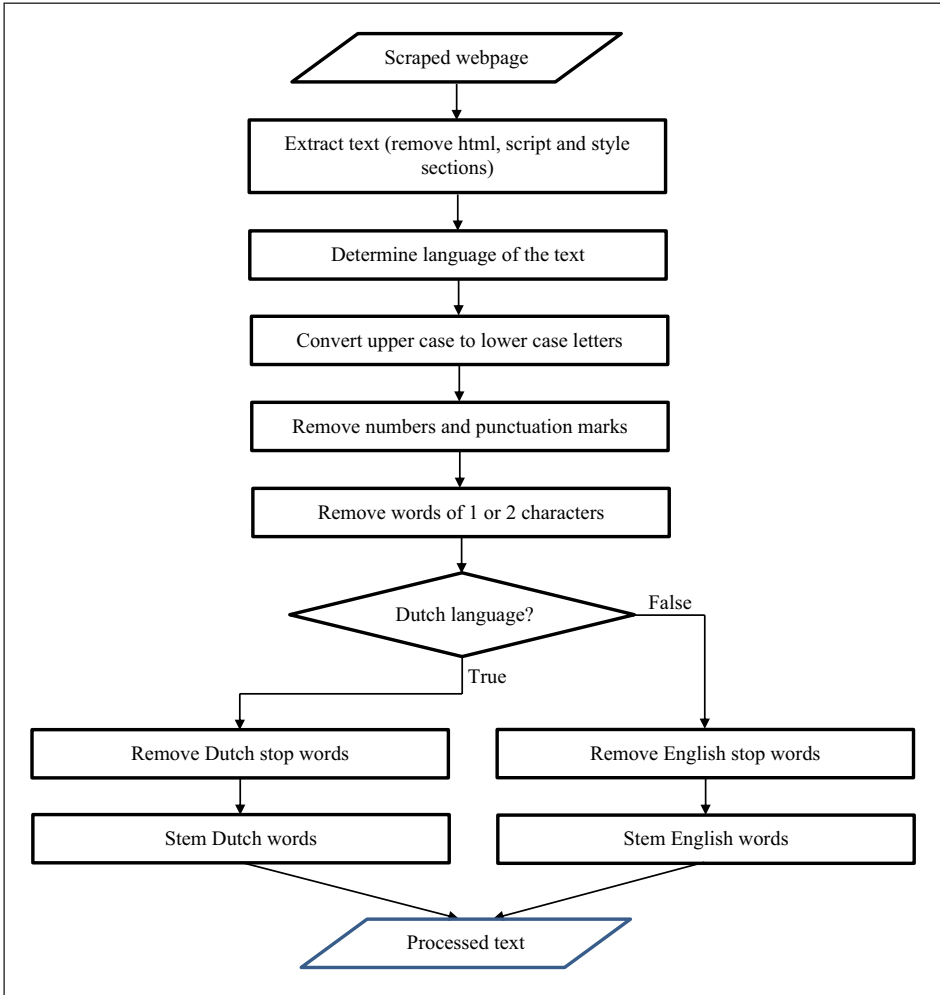


Figure 1. Flow schema of the processing of the texts extracted from scraped webpages. Because a large majority of the non-Dutch texts were found to be written in English, these texts were all processed as such.

statistics and have been studying online platforms for at least two years. To identify the platform and non-platform organizations, the experts review information on websites for a set of organizations. To identify online platforms, the experts essentially answer the two questions included in the Appendix. Websites for which the experts’ opinions are not identical (around 15%) are discussed by all experts after which a joint decision is made. Based on the definition of an online platform (see the introduction) this leads to a set of

$$Platforms = \{DPlatform_j = 1 | Text_j; j = 1, \dots, N\}$$

where the 0 or 1 variable $DPlatform$ has the value of 1 if the organization is characterized as an online platform according to the judgment of the experts. It is based on the multidimensional variable $Text$, which is composed of the combination of words extracted from the website of the organization. The elements in this set are referred to by the subscript j . In total there are N platforms in the set $Platforms$. The experts also assemble a second set of size N , named $Non_Platforms$, for which the 0 or 1 variable $DPlatform$ has the value of 0. In the combined data set, the sets $Platforms$ and $Non_Platforms$ have an equal number of elements and they do not overlap. The combined data set, of size $2N$, created by experts consists of

$$Combined\ data = \{Platforms, Non_Platforms\} \quad (1)$$

In the second step of text processing, a supervised generative model-based approach (Gentzkow et al. 2019) is applied to a large random sample of the combined data set; this is 80% in our case. The sample is referred to as the training data. The model-based approach reduces the multidimensional variable $Text$ to a lower-dimensional variable Z (see next paragraph), such that the new variable Z discriminates the elements of platform versus non-platform organizations. There is a whole range of machine learning algorithms available that can be applied to computational efficiently perform this task (Pedregosa et al. 2011). As a result, the vector Z consists of variables—words—that characterize the elements that belong to the set $Platforms$ versus those that are part of the set $Non_Platforms$. In addition, there usually is a vector of estimated weights $\hat{\theta}$, obtained through a machine learning algorithm, which are used to predict the dichotomy $DPlatform_i = 1$ versus $DPlatform_i = 0$; see Gentzkow et al. (2019) for more details. It leads to the probability that an organization is an online platform, conditional on Z_j and $\hat{\theta}$

$$P_j = Prob(DPlatform_j = 1 | Z_j, \hat{\theta}) \quad j \in \{Platforms, Non_Platforms\} \quad (2)$$

As is usual in machine learning, an independent sample, referred to as a test set (also known as a holdout set), is used to check the performance of Equation (2). Here, the test set is the 20% part of the combined data that remained after selecting the training data.

In the third step, for the entire population of organizations, the statistical model of Equation (2) is used to predict the probability of $DPlatform_i = 1$ given Z_j and $\hat{\theta}$.

$$\hat{P}_i = Prob(DPlatform_i = 1 | Z_i, \hat{\theta}) \quad i \in \{population\ of\ organizations\} \quad (3)$$

The elements of the population are referred to by the subscript i . Depending on the machine learning algorithm used, either a binary value or a value in the range 0 to 1 is produced. In the latter case, all organizations for which the estimated probability \hat{P}_i is above a specifically defined threshold are classified as online platform organizations. Usually, a threshold value of 0.5 is used for this purpose, but this is not always the best

choice; for instance when dealing with highly imbalanced data, such as data with only a limited number of positive cases (Kuhn and Johnson 2013, chap. 16).

Following the three steps described above, there are two major outcomes regarding the identification of online platform organizations. First, there is the set of platform organizations, included in the test set, that have been classified by the statistical model (Equation 2). This set solely consists of organizations for which the correct classification is known as these have been assessed by the experts. This result is used to determine the performance of the model developed, which, in the case of accuracy, refers to the correct identification of platform and non-platform businesses in the total of all businesses classified. Here, one wants to obtain a model with the highest accuracy possible. Second, there is the outcome of the classification of the organizations in the population that were not assessed by the experts. Some of them are identified as online platform organizations (Equation 3). How well the classification model performs on the unobserved organizations, for example, the second outcome, affects the findings tremendously. Especially for the latter results (Equation 3), there is usually no information on the size of the type-I and type-II errors. The type-I errors, that is, the false positives, consist of businesses that are identified by the model as online platforms, even though they are non-platform organizations. The type-II errors, that is, the false negatives, are organizations not classified as a platform by the model but they are actually online platforms. Given that the number of platform organizations is expected to be relatively limited (Heerschap et al. 2021), it will be hard to assess the type-II errors through a survey.

It would be possible, however, to estimate the size of the type-I error. In order to do so, we sent all organizations, after extensive checking, with a value above the threshold (obtained by Equation 3) the Dutch Online Platform questionnaire. Businesses themselves will disclose in the survey whether they can be categorized as an online platform organization, according to the OECD definition. We used the first two questions in this survey for this purpose. The questions are included in the Appendix. Thus there are two outcomes of the external validation. Either a business confirms the outcome of the text-based classification, “the model is right”, or it reports that it is not a platform organization, which leads to a type-I error. More formally, the external validation is based on the external measure \hat{P}_i of Equation (3) which is confronted with externally collected information on the latent 0 or 1 variable $Platform^*$ for the organizations above the threshold. This variable gets a value of 1 if the business confirms that it can be characterized as an online platform. Here, only the information is used of businesses that: (A) indicate they are a platform and (B) are not included in the combined data set created by experts.

Overall, the external validation of the businesses above the threshold provides two important pieces of information that will be examined in the empirical analysis. First, it gives an estimate of the fraction of type-I errors in the estimated platform organizations obtained after applying the model. Second, it leads to an estimate of the statistical association between \hat{P} of Equation (3) and the latent variable $Platform^*$. There will be an indication that the text-based classification leads to satisfactory estimation results if there is a positive association between \hat{P} and $Platform^*$. We are aware that there are other potential sources of bias in the classification process studied, but the main focus of the study described in this article is on identifying online platforms and the effect of type-I errors in particular.

3. Results

3.1. Analysis of Text-Based Results

3.1.1. Step 1: Combined Data Set Creation and Data Collection

In the first step, we constructed a data set containing known examples of platform and non-platform organizations (Equation 1). The websites assigned to the organizations in the Business Register of Statistics Netherlands were an important starting point. Based on this register, the content of the associated websites, the findings of some initial studies (Heerschap et al. 2021), and their expertise, three business statistics employees of Statistics Netherlands created a set of 590 online platform organizations and identified 303 non-platform organizations, with very similar characteristics, during this process. To the latter, a random sample of 287 non-platform organizations, from the websites linked to the Business Register, were additionally added. The websites in this sample were manually checked to assure they were active, were not already included, and were definitely of non-platform organizations. This resulted in a combined data set of in total 1,180 organizations with a website. The organizations and websites included in the combined data set were removed from the large Business Register linked data set in subsequent analysis.

The combined data set of 1,180 organizations was used as the training and test set to develop an online platform text-based classifier. To obtain the texts required for our study, the websites in the combined data set, for example, all 1,180, and all websites linked to the Business Registers units, a total of 960,588, were attempted to be scraped. We found that all websites in the combined data set and 629,284 (66%) of the websites linked to the Business Register could be scraped. The websites that could not be scraped were found to be no longer active.

3.1.2. Step 2: Text Processing and Model Development

In the second step, the variable Text is reduced to the variable Z of lower dimension, and the estimate $\hat{\theta}$ is obtained (Equation 2). This step starts with extracting the texts from the webpages collected in the previous step. For the combined data set, the text in a total of 1,138 websites (96%) could be extracted and processed. The 1,138 texts contained 569 (50%) platform (positive) and 569 (50%) non-platform (negative) cases. Previous studies have shown that a data set constructed in this way is very well suited to determine if a particular text-based classification method is able to differentiate between the texts in the positive and negative cases for the topic studied (Daas and van der Doef 2020; Kuhn and Johnson 2013). Of the 1,138 texts, an 80% random sample was used for model development; this is the so-called training set and included 910 cases. The remaining 20% were used as the test set (see below). Model development required the creation of a document-term matrix in which the rows corresponded to the organizations webpages and the columns to all of the unique words in the training set; see Materials and Methods for more details. Words that occurred fewer than fifty times in the training data were removed. The resulting document-term matrix had a dimension of 910 rows \times 570 columns and 910 rows \times 4,300 columns for those based on the text of the main page only

and those based on the texts extracted from all pages collected on the website, respectively. Subsequently, a whole range of different machine learning methods were trained to discern platform from non-platform websites in the best possible way. The performance of those models was evaluated on the unseen 20% test set of 228 cases. Accuracy, for example, the number of correctly classified cases of the total number of cases included, Precision, for example, the number of positives correctly classified of the total number of cases classified as positives, and Recall, for example, the number of positives correctly classified of the total number of positive cases included, were used as the most important evaluation metrics.

The metrics for various trained classification methods, such as Naive Bayes, Logistic Regression, Support Vector Machines, Regression Trees, and Neural Networks, were compared. During this comparison, the effect of various processing steps on the texts and the choice to use only the words on the main page or the words on all pages scraped for a website were compared. Hyperparameter tuning, via a Tree of Parzen Estimators (Bergstra et al. 2011) followed by five-fold cross-validation, was used to assure the best possible outcome was obtained for each method. It was found that a trained Support Vector Machine (SVM) model with a linear kernel produced the best results when: (A) the words on all pages collected from of a website were used, (B) the words were stemmed, and (C) only words of three or more characters were included. Hyperparameter optimization revealed that the standard settings for this method already resulted in the best performance. For the SVM model, an accuracy of 82% ($\pm 2\%$) was obtained on the test set (as defined above in this section). The standard deviation was determined by repeating the entire procedure on resamples, with replacement, one thousand times. The precision was 84% ($\pm 3\%$), and the recall was 79% ($\pm 4\%$). Applying deeply trained Bidirectional Encoder Representations from Transformers (BERT) or its Dutch version BERTje (Fialho et al. 2020) did not produce better results. The creation of two language-specific models, one for Dutch and one for English websites, did also not improve the overall findings. Even though the SVM approach did not work perfectly, it produced the best statistical model to identify online platform websites of all options and combinations tested. Including Word Embeddings derived features did not improve the classification findings even after additional hyperparameter optimization. The SVM model obtained provided a score of being an online platform website that was scaled to a probability via a five-fold cross-validation procedure (Platt 2000). That probability (Equation 2) is, from here on, indicated as "PWebsite". This is a value between 0 and 1 which can be converted to a class label by using a threshold. Any value above that threshold is considered a positive (platform) case.

The findings of the SVM model were additionally checked by studying the distribution of the probabilities on the test set. This revealed a, somewhat noisy, U-shape indicating that the two cases could be separated fairly well. In addition, the ten words with the highest positive and highest negative coefficients used by the SVM model were inspected (see Table 1). The findings for the words with high positive coefficients indicate that the trained model picked up the intended classification topic. The words with high negative coefficients are indicative of a heterogeneous group of websites which is not an unexpected finding as there is a whole range of non-platform websites. Let it be clear that it is the combination of words, that remain after processing website texts, that is used by

Table 1. Words with the Ten Highest and Ten Lowest Coefficients in the Trained Support Vector Machine Model.

Word	Coefficient	Word	Coefficient
platform	2.526	portfolio	-1.202
account	1.690	phone	-1.113
help	1.653	info	-1.098
crowdfunding	1.587	skip	-1.047
register	1.551	approach	-1.029
login	1.236	year	-1.025
entrepreneur	1.175	wordpress	-0.982
deal	1.152	since	-0.981
ask	1.149	p.o. box	-0.964
neighborhood/vicinity	1.143	customization	-0.959

the SVM model to produce a probability score. As a consequence, the presence of a login module on a particular website alone does not have to result in a score high enough to be considered an online platform, for instance when other important words (such as platform) are absent.

3.1.3. Step 3: Prediction

In the third step, the statistical model is used to predict which organizations in the population are online platforms (Equation 3). Hence, the SVM model was applied to the texts extracted and processed from the huge set of web pages linked to the Business Register, while excluding those included in the combined data set. From the 629,284 websites linked a total of 10,964,998 pages could be scraped; indicating an average of slightly more than seventeen pages collected per website. The web pages were processed according to the optimal procedure described above which resulted in a text file for each website. Of those files, 594,574 (94.5%) contained ten words or more. These files were classified with the SVM model developed and, for each case, the probability of being an online platform website was determined. This resulted in 41,811 (7%) websites with a value above the standard cut-off value of 0.5 which is usually used as a threshold to indicate the positive cases. However, the histogram of the distribution of the probabilities for these websites indicated a strongly negatively skewed distribution (Figure 2). Because of this finding and the fact that the model was developed on an equal number of platform and non-platform websites, a situation that is expected to be far off from the platform-non-platform ratio of websites linked to the Business Register, made clear that the classification findings needed to be studied in more detail. This was done by drawing random samples of fifty websites in nine probability ranges, each 0.1 wide, and manually inspecting the websites selected. This revealed that in these samples, online platform websites started to occur at probabilities values of 0.8 and higher. In the lower value ranges, none were detected in the samples drawn. From this, it is clear that the model obviously overestimates the number of platform websites in the Business Register when a value of 0.5 is used as the positive detection threshold. This is an important finding as it reveals that

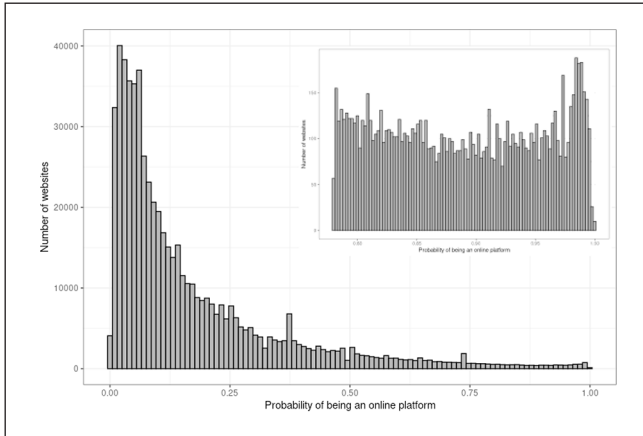


Figure 2. Histogram of the model-based probabilities of being an online platform website.

Note. Information is used for the 594,574 websites with ten words or more linked to the Business Register. The insert shows the findings for websites with a probability above 0.78 and reveal a peak around 0.98 to 0.99.

the model behaves differently on the Business Register data compared to the training and test data. This not only has interesting research applications, described in more detail in Puts and Daas (2021), but also suggests applying a higher threshold value for online platform identification. The fact that the lowest point in the probability distribution is somewhere located around 0.90 corroborates this observation (Figure 2). We found that the number of websites with a \hat{P} value above 0.8 is 9,129 (1.5%). In addition, the probability distribution reveals a small peak around 0.99 (insert in Figure 2), which suggests the occurrence of a group with very high \hat{P} results.

3.2. External Validity

3.2.1. Survey and Selection Procedure Used

To determine the extremal validity of the findings described above, an additional (new) source of information is needed. For this, we use the response of the organizations to the Dutch Online Platform survey of Statistics Netherlands (Klijns 2021). First, we report the selection procedure followed for the organizations given in Table 2. We start with the 9,129 websites identified as those of an online platform organization by the statistical model, all of which have a fitted probability \hat{P} above a value of 0.8. Three subsequent selection steps were taken to construct the final survey population (see Table 2). First, websites with adult content were removed because these were not considered to belong to the target population. The URLs of these websites were checked for the occurrence of words or parts of words typical for adult content websites. Any sites that contained one (or more) of these words or word parts were removed. Next, the relation of the website with the, to be surveyed, legal units in the Business Register was meticulously checked. Because websites were assigned at a very detailed level, that is, the local unit level, for a

Table 2. The Selection Procedure Followed.

	Number	Percentage selected	Percentage surveyed
Websites with a platform probability of at least 0.8	9,129	100.0	
Removal of adult content websites	7,764	85.0	
Distinct legal entities	6,057	66.3	
Legal entities approached with the platform survey	4,385		100.0
Response from the legal entities	2,997		68.3
Usable response for empirical analysis	2,708		61.8
Usable response excluding websites in the “combined data set”	2,631		60.0

Note. Only websites with a platform probability of 0.8 or higher, as indicated by the model developed, were included.

considerable number of (legal) units, that is, 1,707 organizations as indicated in Table 2, multiple websites were found. For these businesses only one website was randomly selected to avoid sending multiple surveys to one business. In the final step, it was assessed whether information about the business unit to which the legal unit belonged could be retrieved from the business register. If this was not the case, the legal unit was excluded from the survey. In the end, a total of 4,385 organizations were approached to participate in the Dutch Online Platform survey. The response to the survey was 68% (2,997 organizations), which is relatively high for a business survey. Of all responding organizations, 289 respondents were excluded from further analysis since information on any of the essential variables required in the subsequent analysis procedures was missing in their response. This resulted in a selection of 2,708 organizations to be used in our subsequent statistical analysis.

3.2.2. Type-I Errors

To validate the results of the text-based classification, we examine the relative size of the type-I errors, which are the false positives of observing an online platform organization. There are two reasons for a type-I error. First, in the first step of the estimation procedure, the expert opinion leads to a wrong assessment of some of the organizations that are considered platform organizations. Second, in the third step of the estimation procedure, in which the empirical model is applied to the entire population, some of the organizations are mistakenly predicted as platform organizations.

To inspect the false positives, the responses to two questions included in the Dutch Online Platform survey are used; see Appendix. First, we report on the percentage of false positives. Out of the 2,708 organizations with a usable response, 2,064 organizations (76.2%) responded negatively to the question of whether their website is an online platform. This result suggests a substantial percentage of (potential) false positives. So a total of 644 platform organizations were initially found.

Next, a comparison is made between the combined data (used for model development in steps 1 and 2) and the classified organizations. The latter data are retrieved by step 3 of the statistical procedure, and in what follows they are referred to as the “predicted data.” Of the 2,708 responding organizations, there are seventy-seven organizations included in the combined data set of step 1 and 2. This means there are 2,631 responding organizations in the predicted data set. The percentage of organizations that indicated they are not an online platform is 54.5% for the organizations in the combined data and 76.9% for the organizations in the predicted data; a total of 2,022. This brings us to the second result, that organizations in the combined data, obtained through expert assessment in step 1, better reflect the platforms than the remaining (“predicted”) organizations obtained by prediction in step 3.

Next, we focus on the relationship between the false positives and \hat{P} . To prevent any potential biases by the expert assessment (step 1), we restrict ourselves to the 2,631 responding organizations in the predicted data set. Hence, only information on the 2,631 organizations is used in Tables 3 and 4. Table 3 gives an overview of the findings of the predicted data broken into several categories of \hat{P} . The average for *Platforms* seems to be positively related to the probability of the website being a platform, as indicated by the model, ranging from an average of 15.9% for the category 0.80–0.839 to 36.7% for the category 0.96 and higher. This brings us to the third result, namely that there seems to be a smaller fraction of false positives for organizations with a larger \hat{P} value. This relationship will be explored in more detail in the next section. Table 3 also reports on some additional characteristics of the organizations included. A relatively large part of the organizations is comparatively small. About 25% of the organizations have one employee only; 14% of the organizations are large and have at least fifty employees.

Finally, we reassessed a sample of the false positive organizations in the predicted data set by re-evaluating their website by the experts of Statistics Netherlands; a so-called second opinion. A random sample of one hundred (false positive) organizations was drawn from the 2,022 organizations identified as false positive in the predicted data set. The sample was stratified by the \hat{P} value categories used in Table 3. Table 4 reports the percentage positives for the false positives in each category according to the combined opinion of the experts. As a fourth result, it was found that, for some of the organizations, there is a disagreement between the opinion of the experts and the organization itself on whether their website is an online platform. According to the experts’ reassessment, a weighted average of about 21% of the false positives in the predicted data set is characterized as a positive (a platform); indicating that it is a “false false positive”; for example, an actual platform organization. For organizations with a \hat{P} value close to one, the disagreement is most strong. In the highest \hat{P} value category, 65% of the false positives are online platforms according to the experts’ reassessment (Table 4).

Based on these empirical findings, we computed a lower and an upper bound for the percentage of false positives. There are two extreme situations. If we assume that the response to the survey gives a correct representation of whether or not the organizations surveyed are online platforms, the percentage of false positives is 76.9%. Alternatively, if we interpret the expert opinion as the gold standard, the percentage of false positives is only 60.8%; ($76.9 \times (100 - 21) / 100$). Based on the information available, it is expected that the actual percentage of false positives is somewhere between 60.8% and 76.9%. This

Table 3. Summary Statistics.

	Number of organizations	Percentage platform
Total	2,631	23.1
Platform probability		
0.800–0.839	753	15.9
0.840–0.879	649	20.8
0.880–0.919	521	25.7
0.920–0.959	411	27.0
≥ 0.960	297	36.7
Number of employees		
≤ 1	1,366	24.1
1.1–4.9	594	24.9
5.0–19.9	328	22.6
20.0–49.9	144	23.6
≥ 50	199	12.1
Branch		
Wholesale and retail trade; repair of motor vehicles	420	12.1
Information and communication	745	31.7
Renting and trading real estate	189	15.9
Consultancy, Research and Other Specialist Business Services	444	24.8
Rental of movable property and other business services	166	34.9
Other branches	667	18.6
Legal form		
Proprietorship	971	20.9
Private company	1,192	24.0
General partnership	267	27.0
Other legal form	201	23.9

Table 4. Experts' Second Opinion on the False Positives.

	Number (%) no platform according to survey response	Percentage positive according to experts' check ^a
Platform probability		
0.800–0.839	633 (84.1)	15.0
0.840–0.879	514 (79.2)	10.0
0.880–0.919	387 (74.3)	25.0
0.920–0.959	300 (73.0)	20.0
≥ 0.960	188 (63.3)	65.0
Total	2,022 (76.9)	21.0 ^b

^aFor a random sample of one hundred organizations (twenty in each stratum of platform probability).^bWeighed by number of organizations in the five strata.

suggests that the population contains at least 1,121 to 1,455 online platform organizations after correcting for false positives. Since there is no information on the false negatives, we emphasize that this value is, very likely, an underestimation of the true number of platform organizations. The group of organizations not surveyed, that is, those with a \hat{P} below 0.8, could potentially still include some online platform organizations. However, because of the expected low numbers in this group, a costly and time-consuming survey is needed to reliably obtain information on the number of false negatives; that is, the type-II errors among those organizations.

3.2.3. Statistical Association

In this subsection, we measure the statistical association between \hat{P} of Equation (3), which was obtained through the empirical analysis of the text-based classification findings, and the latent variable $Platform^*$, obtained by the survey. The requirement to receive a questionnaire was that the \hat{P} value of the text-based classification was at least 0.8 so the fitted probability value \hat{P} has a range of 0.8 to 1.0. There is an indication that the text-based analysis leads to satisfactory estimation results if there is a positive association between \hat{P} and $Platform^*$. For this, only the organizations in the predicted data set were used to ensure there was no contamination of the expert opinion from the first step of the empirical analysis.

The starting point for the empirical specification is the response of the organizations to the survey questions about whether their website is an online platform. The dependent variable is a 0 or 1 indicator that gets the value of 1 if there is a positive response to the two questions shown in the Appendix. If that is not the case, the organization is characterized as a false positive. The variable regressed on is the fitted probability value of the platform organization, reported in subsection 3.1.3, as well as some control variables. The average of the dependent variable is 0.231 (Table 2), which makes it sufficiently large to specify the regression as a linear probability model (LPM)

$$Y_platform_i = \beta_0 + \beta_1 \hat{P}_i + \gamma' X_i + u_i \quad i = 1, \dots, M \quad (4)$$

where subscript i refers to the i -th organization; there are M organizations. $Y_platform$ is a 0 or 1 indicator variable which is one, if the organization responds in the survey that it performs online platform activities. \hat{P} is the fitted probability value that was obtained through the text-based classification; it may have a slight attenuation bias due to the uncertainty of \hat{P} . The vector X contains variables from the Business Register: firm size (four categories), economic sector (three categories), and legal status (two categories). u is an idiosyncratic error term. The advantage of an LPM, which is a linear regression equation for which the dependent variable is a 0 or 1 binary variable, is that it is particularly useful in specifications with discrete explanatory variables to detect a non-linear relationship. To apply a heteroskedasticity robust inference, we calculate the robust standard errors of the estimated regression parameters, for which we will report the White standard errors (Wooldridge 2010).

Although a valid probabilistic measure of \hat{P} implies that the estimated β_1 has a positive sign, we can be more specific about its value. There is a strong indication of a valid empirical model for the platform probability since the change of \hat{P} from 0.8 to 1.0 is associated with a change of the average of $Y_platform$ for all organizations by 0.2 percentage points. In other words, there is an indication of a one-to-one relationship between $PWebsite$ and $Y_platform$. This hypothesis will be tested below.

Next, we discuss the parameter estimates of Equation (4) that are reported in Table 5. According to the various specifications, the marginal effect of the estimated parameter on the platform probability ranges from 0.011 to 0.013. It means that an increase of the probability \hat{P} by 1 percentage point corresponds to an increase of the probability of being a platform by 1.1 to 1.3 percentage points. The estimated parameters presented in Table 6 confirm this outcome. The platform probability is distinguished into five categories. The estimates indicate that relative to the reference category of 0.8 to 0.834, the difference for the upper category is 19.5 percentage points. Thus, for a change of \hat{P} from 0.8 to 1.0, there is also an increase of a positive survey response by about 20 percentage points. Remarkably, the categories 0.88 to 0.919 and 0.92 to 0.959 give similar parameter estimates (about a 10% difference relative to the reference category). As a robustness check, using a Logit specification for Equation (4) leads to the same estimated marginal effects.

4. Discussion

In this article we have reported an ex-post empirical analysis of the validation of a sub-population that was identified through the combination of web scraping and text-based classification. Such a statistical procedure is a useful tool in case it is hard to identify the target population of research when conventional sampling methods from a predefined population, such as stratified sampling or cluster sampling, cannot be applied. Here, it helps that the procedure can be easily applied to large amounts of data. Indeed, for the application that is described in this article, it became clear that online platform activities are unevenly distributed economy-wide.

The estimates of the text-based classification procedure on online platform organizations lead to satisfactory outcomes after a number of additional selection steps. We consider the most relevant and least relevant words identified by the statistical model to be plausible. Furthermore, the distribution of fitted probabilities of being an online platform for the population of businesses gives a bimodal distribution. Our analysis also confirmed the earlier observation that a model developed on 50% positive (platform) and 50% negative (non-platform) cases, behaves differently when the model is applied to real-world ratios of these cases (Puts and Daas 2021). For online platform detection, which occurs much less than 50% in our “real-world” data, this results in an overestimation of the number of positive cases and makes the study, and consequently the reduction, of the number of false positive cases an important topic of the work described in this article. Here, it becomes clear that the model is able to identify a rare occurring group of (potential) online platform organizations in a very large population.

Table 5. Estimates Equation (4), Part I.

	Estimate (robust standard error)		
	Model 1	Model 2	Model 3
Intercept	-0.823*** (0.133)	-0.744*** (0.135)	-0.655 (0.512)
Platform probability	1.196*** (0.152)	1.073*** (0.150)	0.973 ^a (0.585)
Number of employees			
≤ 1		0.025 (0.026)	-0.088 (0.407)
5–19.9		-0.017 (0.030)	-0.135 (0.502)
20–49.9		-0.03 (0.041)	-1.05 ^a (0.623)
≥ 50		-0.128*** (0.032)	1.026 ^a (0.542)
Branch			
Wholesale and retail trade; repair of motor vehicles		-0.104*** (0.026)	0.297 (0.438)
Information and communication		0.059* (0.026)	0.244 (0.422)
Renting and trading real estate		-0.093** (0.034)	0.205 (0.582)
Rental of movable property and other business services		0.107* (0.042)	0.235 (0.627)
Other branches		-0.065* (0.026)	0.682 ^a (0.404)
Legal form			
General partnership		0.086* (0.036)	-0.097 (0.592)
Private company		0.064** (0.023)	-0.459 (0.359)
Other legal form		0.111** (0.036)	-1.595* (0.641)
Interactions platform probability × employees			
PP × one or less			0.130 (0.465)
PP × more than five, less than twenty			0.137 (0.577)
PP × more than twenty, less than fifty			1.164 (0.722)
PP × more than fifty			-1.320* (0.619)
Interaction platform probability × Branch			
PP × Wholesale and retail trade; repair of motor vehicles			-0.458 (0.507)
PP × Information and communication			-0.211 (0.482)
PP × Renting and trading real estate			-0.339 (0.669)
PP × Rental of movable property and other business services			-0.140 (0.712)
PP × Other branches			-0.847 ^a (0.462)

(Continued)

Table 5. (Continued)

	Estimate (robust standard error)		
	Model 1	Model 2	Model 3
Interactions platform probability \times Legal form			
PP \times General partnership			0.205 (0.677)
PP \times Private company			0.590 (0.409)
PP \times Other legal form			1.948** (0.739)
Goodness-of-fit			
R-squared	0.024	0.062	0.069
Percentage predicted platform probability between 0 and 1	100.0	99.1	99.7

Note. Dependent variable is the 0 or 1 variable $Y_{platform}$.

*** $p < 0.001$. ** $p < 0.01$. * $p < 0.05$. ^a $p < 0.1$.

Our results show that applying the classification model seriously reduces the initial population of around 600,000 organizations to a set of a bit more than nine thousand organizations; a more than 60-fold reduction. The latter data set is highly enriched in platform organizations and could be, after some additional checking and selection, almost completely surveyed. To statisticians that want to apply our empirical approach, we strongly advise first making sure that the machine learning model is able to discern between the positive and negative cases of the topic studied in the best possible way. This requires not only a data set with typical positive and negative examples but also clear negative examples that, at first sight, resemble the positive cases reasonably well. In this way, one tries to make sure that only the relevant (and hence important) words are to be included in the model. When such a model is applied to the entire population, one subsequently needs to carefully check the external validity of the model by, for instance, manually inspecting websites.

We list a number of methodological learning points from our study. First, the quality of the data set used to build a classification model is important to get a bimodal distribution of predicted outcomes. Second, the fraction of units in the target population must be sufficiently large to enable adequate detection. Third, the external validation indicates a positive association between the fitted probability of the text-based model and a positive response to the survey question on being an online platform organization. Fourth, the fraction of false positives is large. Fifth, we observe a remarkable disagreement between the organization and the experts' opinion on the question of whether the organization can be characterized as an online platform. The latter could, for instance, be caused by the fact that the first two questions in the Dutch Online Platform survey focus on deriving if the organization can be characterized as an online platform. Any respondents that do not want to answer the remaining questions of the survey can simply end the survey by stating the business is not a platform organization. Future research will focus on dealing with the last three learning points since we expect it will substantially enhance the application

Table 6. Estimates Equation (4), Part 2.

	Estimate (robust standard error)	
	Model 1	Model 2
Intercept	0.159*** (0.013)	0.137*** (0.034)
Platform probability		
0.840–0.879	0.049* (0.020)	0.042* (0.021)
0.880–0.919	0.098*** (0.023)	0.096*** (0.023)
0.920–0.959	0.111*** (0.026)	0.092*** (0.025)
≥ 0.960	0.208*** (0.031)	0.187*** (0.030)
Number of employees		
≤ 1		0.024 (0.026)
5–19.9		–0.018 (0.030)
20–49.9		–0.034 (0.041)
≥ 50		–0.129*** (0.032)
Branch		
Wholesale and retail trade; repair of motor vehicles		–0.104*** (0.026)
Information and communication		0.062* (0.027)
Renting and trading real estate		–0.091** (0.034)
Rental of movable property and other business services		0.111** (0.042)
Other branches		–0.063* (0.026)
Legal form		
General partnership		0.085* (0.036)
Private company		0.063*** (0.023)
Other legal form		0.110** (0.036)
Goodness of fit		
R-squared	0.023	0.061

Note. Dependent variable is the 0 or 1 variable $Y_{platform}$.

*** $p < 0.001$. ** $p < 0.01$. * $p < 0.05$.

of machine learning in official statistics and the study of rare business subpopulations. Finally, we leave to further empirical research the size of type-II errors.

Authors' Note

The views expressed in this article are those of the authors and do not necessarily reflect the policies of Statistics Netherlands. Part of this article is based on the Dutch Online Platform survey carried out on behalf of the Dutch Ministry of Economic Affairs and Climate Policy.

Acknowledgements

The authors like to thank Kelby Çakim for her excellent assistance during the beginning of the project, Tim de Jong for his assistance in the hyperparameter optimization part of the study, and Marco Puts for stimulating discussions. Comments by Joep Burger, Javier Garcia-Bernardo, Joram

Vuik, Ger Snijkers, and Yvonne Gootzen are gratefully acknowledged. We highly appreciate the comments and suggestions of the anonymous reviewers and (associate) editors on earlier versions of the article which have greatly improved the final manuscript.

Funding

The authors received no financial support for the research, authorship, and/or publication of this article.

References

- Aggarwal, C. C. 2016. "Mining Text Data." In *Data Mining: The Textbook*, edited by C. C. Aggarwal, 429–55. New York: Springer.
- Allen, C., and T. Hospedales. 2019. "Analogies Explained: Towards Understanding Word Embeddings." Proceedings of the 36th International Conference on Machine Learning: ICML, Long Beach, CA, USA, June 11–13. <http://proceedings.mlr.press/v97/allen19a/allen19a.pdf> (accessed November 2023).
- Becue, M., B. Fridlund, A. Fyhrlund, A. Prat, and B. Sundgren. 2004. "Text Mining in Official Statistics." In *Text Mining and Its Applications: Results of the NEMIS Launch Conference*, edited by S. Sirmakessis, 189–204. Berlin: Springer.
- Berg, J., M. Furrer, E. Harmon, U. Rani, and M. Silberman. 2018. *Digital Labour Platforms and the Future of Work. Towards Decent Work in the Online World*. Geneva: International Labour Organization. https://www.ilo.org/wcmsp5/groups/public/—dgreports/—dcomm/—publ/documents/publication/wcms_645337.pdf (accessed November 2023).
- Bergstra, J., R. Bardenet, Y. Bengio, and B. Kégl. 2011. "Algorithms for Hyper-Parameter Optimization." In *Advances in Neural Information Processing Systems 24*, edited by J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger. New York: Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf> (accessed November 2023).
- Cui, L., Y. Hou, Y. Liu, and L. Zhang. 2020. "Text Mining to Explore the Influencing Factors of Sharing Economy Driven Digital Platforms to Promote Social and Economic Development." *Information Technology for Development 27*: 779–801. DOI: <https://doi.org/10.1080/02681102.2020.1815636>.
- Daas, P. J. H., M. J. Puts, B. Buelens, and P. A. M. van den Hurk. 2015. "Big Data and Official Statistics." *Journal of Official Statistics 31*: 249–62. DOI: <https://doi.org/10.1515/jos-2015-0016>.
- Daas, P. J. H., and S. van der Doef. 2020. "Detecting Innovative Companies via Their Website." *Statistical Journal of IAOS 36*: 1239. DOI: <https://doi.org/10.3233/SJI-200627>.
- De Groen, W. P., Z. Kilhoffer, K. Lenaerts, and N. Salez. 2017. "The Impact of the Platform Economy on Job Creation." *Intereconomics 52*: 345–351. DOI: <https://doi.org/10.1007/s10272-017-0702-7>.
- Ducci, F. 2020. *Natural Monopolies in Digital Platform Markets*. Cambridge: Cambridge University Press.
- Eurostat. 2008. *NACE Rev. 2-Statistical classification of economic activities in the European Community*. Available at: <https://ec.europa.eu/eurostat/web/products-manuals-and-guidelines/-/ks-ra-07-015> (accessed March 2024).
- Fialho, P., L. Coheur, and P. Quaresma. 2020. "To BERT or Not to BERT Dealing with Possible BERT Failures in an Entailment Task." In *Information Processing and Management of Uncertainty in Knowledge-Based Systems, Computer and Information Science Book 1237*, edited by M.-J. Lesot, S. Vieira, M. Z. Reformat, J. P. Carvalho, A. Wilbik, B. Bouchon-Meunier, and R. R. Yager, 734–47. Cham: Springer International Publishing.

- García Lozano, M., J. Brynielsson, U. Franke, M. Rosell, E. Tjörnhammar, S. Varga, and V. Vlassov. 2020. "Veracity Assessment of Online Data," *Decision Support Systems* 129: 113132. DOI: <https://doi.org/10.1016/j.dss.2019.113132>.
- Gentzkow, M., B. Kelly, and M. Taddy. 2019. "Text as Data." *Journal of Economic Literature* 57: 535–74. DOI: <https://doi.org/10.1257/jel.20181020>.
- Heerschap, N., B. Klijs, A. Mares, M. van Rossum, and J. Vuik. 2021. "Getting a Grip on the Platform Economy in the Netherlands." The 36th IARIW Virtual General Conference, Online conference, August 23–27. https://iariw.org/wp-content/uploads/2021/07/vanRossum_Paper.pdf (accessed November 2023).
- Howcroft, D., and B. Bergvall-Kåreborn. 2019. "A Typology of Crowdwork Platforms." *Work, Employment and Society* 33: 21–38. DOI: <https://doi.org/10.1177%2F0950017018760136>.
- Klijs, B. 2021. *Monitor Online Platforms 2020 (in Dutch)*. The Hague/Heerlen: Statistics Netherlands. <https://www.cbs.nl/nl-nl/longread/rapportages/2021/monitor-online-platformen-2020> (accessed November 2023).
- Kuhn, M., and K. Johnson. 2013. *Applied Predictive Modeling*. New York: Springer.
- Luiten, A., J. Hox, and E. de Leeuw. 2022. "Survey Nonresponse Trends and Fieldwork Effort in the 21st Century: Results of an International Study Across Countries and Surveys." *Journal of Official Statistics* 36: 469–87. DOI: <https://doi.org/10.2478/jos-2020-0025>.
- OECD. 2019. *Measuring the Digital Transformation: A Roadmap for the Future*. Paris: Organization for Economic Co-operation and Development. <https://web.archive.oecd.org/2020-07-23/559604-roadmap-toward-a-common-framework-for-measuring-the-digital-economy.pdf> (accessed March 2024).
- OECD. 2020. *A Roadmap Toward a Common Framework for Measuring the Digital Economy*. Paris: Organisation for Economic Co-operation and Development. <https://www.oecd.org/sti/roadmap-toward-a-common-framework-for-measuring-the-digital-economy.pdf> (accessed November 2023).
- Oostrom, L. A. N., A. N. Walker, B. Staats, M. Slootbeek-Van Laar, S. Ortega-Azurduy, and B. Rooijakkers. 2016. *Measuring the Internet Economy in The Netherlands: A Big Data Analysis*. The Hague/Heerlen: Statistics Netherlands. https://www.cbs.nl/-/media/_pdf/2016/40/measuring-the-internet-economy.pdf (accessed November 2023).
- Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. "Scikit-Learn: Machine Learning in Python." *Journal of Machine Learning Research* 12: 2825–30. DOI: <https://dl.acm.org/doi/10.5555/1953048.2078195>.
- Platt, J. C. 2000. "Probabilities for Support Vector Machines." In *Advances in Large Margin Classifiers*, edited by A. Smola, P. Bartlett, B. Schölkopf, and D. Schuurmans, 61–74. Cambridge: MIT Press.
- Puts, M. J. H., and P. J. H. Daas. 2021. "Unbiased Estimations Based on Binary Classifiers: A Maximum Likelihood Approach." The 2021 Symposium on Data Science and Statistics, Online conference, June 2–4. <https://arxiv.org/abs/2102.08659> (accessed November 2023).
- Ritzen, J. H. G. 2007. "Statistical Business Register: Content, Place and Role in Economic Statistics." The 3rd International Conference on Establishment Surveys, Montréal, Canada, June 18–21. <https://ww2.amstat.org/meetings/ices/2007/proceedings/ICES2007-000144.PDF> (accessed November 2023).
- Rochet, J.-C., and J. Tirole. 2003. "Platform Competition in Two-Sided Markets." *Journal of the European Economic Association* 1: 990–1029. DOI: <https://doi.org/10.1162/154247603322493212>.
- Snijkers, G., M. Bavdaž, S. Bender, J. Jones, S. MacFeely, J. W. Sakshaug, K. J. Thompson, and A. van Delden. 2023. *Advances in Business Statistics, Methods and Data Collection*. Hoboken, NJ: Wiley.

- Snijkers, G., G. Haraldsen, J. Jones, and D. K. Willimack. 2013. *Designing and Conducting Business Surveys*. Hoboken, NJ: Wiley.
- Sutherland, W., and M. H. Jarrahi. 2018. "The Sharing Economy and Digital Platforms: A Review and Research Agenda." *International Journal of Information Management* 43: 328–41. DOI: <https://doi.org/10.1016/j.ijinfomgt.2018.07.004>.
- Tourangeau, R., B. Edwards, T. P. Johnson, K. M. Wolter, and N. Bates. 2014. *Hard-to-Survey Populations*. Cambridge: Cambridge University Press.
- United Nations. 2019. *Digital Economy Report 2019, Value Creation and Capture: Implications for Developing Countries*. New York: United Nations. https://unctad.org/system/files/official-document/der2019_en.pdf (accessed November 2023).
- United Nations. 2020. *Guidelines on Statistical Business Registers*. Final Draft. New York: United Nations. https://unstats.un.org/unsd/business-stat/SBR/Documents/UN_Guidelines_on_SBR.pdf (accessed November 2023).
- Wooldridge, J. M. 2010. *Econometric Analysis of Cross Section and Panel Data*. Cambridge: The MIT Press.
- Wu, M.-J., K. Zhao, and F. Fils-Aime. 2022. "Response Rates of Online Surveys in Published Research: A Meta-Analysis." *Computers in Human Behavior Reports* 7: 100206. DOI: <https://doi.org/10.1016/j.chbr.2022.100206>.

Received: July 2022

Accepted: November 2023

Appendix

Questionnaire items and definition included in the Dutch Online Platform survey to assess whether the organization's website is an online platform.

An online platform is a website or app where different people, organizations, or companies come into contact with each other and can be linked to each other. Goods, services, or information can then be exchanged via the online platform. The online platform usually does not supply these goods, services, or information itself, but mainly acts as an intermediary.

1. Does your website support or mediate the exchange of goods, services of information between persons, firms, or organizations?
>>This can involve mediation or support in the sale of goods, bringing residents, patients and family together, crowdfunding, dating, new friendships, renting out accommodations, borrowing things etc.<<
 1. Yes
 2. No
2. Are you or is your organization the only provider of the goods, services of information on your website?
 1. Yes
 2. No, in addition to our own supply, there is also supply from other parties
 3. No, there is only offer from other parties

Platform confirmed: question 1 = "Yes," question 2 = "No, in addition.." or "No, there is.."