

Money laundering and the detection of bad companies: a machine learning approach for the risk assessment of opaque ownership structures

Maria Jofre, Antonio Bosisio, Michele Riccardi, Stefano Guastamacchia *

**All authors: Transcrime – Università Cattolica del Sacro Cuore¹*

Beneficial ownership (BO) transparency is a central issue in the anti-money laundering (AML) debate worldwide. In the last 15 years, various measures have been implemented to minimize the opacity of corporate ownership, lastly the establishment of BO registers in various jurisdictions. However, despite all this emphasis, the evidence in support of this concern is limited to few case studies and the empirical research is lacking. The present paper, resulting from the research activity of EU-funded Project DATACROS, addresses this gap. It implements a machine learning approach to, first, develop novel indicators to measure various facets of the opacity of corporate ownership and, second, validate them against evidence of illicit activity. This is done through an empirical analysis of more than 3 million firms registered in nine European countries, and evidence of sanctions and enforcement measures on these firms and/or their owners. Results show that corporate ownership opacity is a relevant factor explaining the likelihood that a firm (or its owners) may be involved in illicit activities. Therefore, the proposed risk indicators have a strong predictive power in identifying companies linked to negative evidence which could be employed in AML/CTF risk assessment activities. The inclusion of macro-level information, such as geographic location and business sector, further improves the understanding of the phenomenon significantly.

Keywords: Money laundering, risk assessment, risk indicators, ownership structure, corporate opacity, machine learning

Draft version prepared for the 2021 Bahamas AML conference – Not to circulate without authors' permission

¹ Authors would like to acknowledge Stefano Barberis and Andrea Maenza for their contribution to the project.

1. Introduction

Corporate ownership transparency has become a central subject in the anti-money laundering (AML) and countering terrorist financing (CTF) global debate since the early 2010s, when, first, some publications (*in primis* the UNODC-World Bank 'StAR' report - van der Does de Willebois et al. 2011) highlighted the role played by shell companies and opaque corporate vehicles in enabling and abetting financial crimes; and, secondly, when the FATF published in 2012 the new (and latest) version of its *Recommendations*, which strengthened the need of having a greater transparency on firms and their beneficial owners (FATF 2012). Since then, a number of measures have been implemented in order to increase the transparency of legal firms, the most important probably being the establishment of beneficial owners (BO) registers in several countries worldwide – and in all European Union (EU) member states following the obligations set by the 4th and the 5th AML Directive (European Parliament and Council of the European Union 2015).

Surprisingly, all the emphasis around the issue, and the subsequent measures taken at global level to address it, are based on poor empirical evidence. Most research literature dealing with corporate opacity and (anti) money laundering relies on a limited number of case-studies (van der Does de Willebois et al. 2011; Transcrime 2018; Steinko 2012). Systematic investigation of the ML/TF risk related to corporate ownership patterns is lacking, apart from some recent works (Garcia-Bernardo et al. 2017; Aziani, Ferwerda, and Riccardi 2020; Riccardi, Milani, and Camerini 2018; Ferwerda and Kleemans 2018). One of the reasons of the lack of empirical research is related to the shortage of reliable corporate ownership data: business registers are national in scope, and do not allow to reconstruct complex cross-border ownership chains; private data providers offer uneven coverage across countries and are often very costly - at least for researchers. The result is that, although we all agree that opaque corporate structures are associated to a higher risk of money laundering, we can count on a poor empirical evidence to prove this statement.

Another interesting pattern in this debate is that the focus has been posed by both regulators and researchers on *who* controls firms, and only very seldom on the techniques, methods and *modi operandi* employed to achieve such control. In other words, all AML/CTF resources have been devoted to the identification and understanding of *who* lies behind a legitimate firm – i.e. the beneficial owner(s) – but not on the *how*. The same BO registers go in this direction, assuming that disclosing the name of the natural person(s) behind a legal entity would be enough to make it fully transparent and minimize the risk that it could be misused for criminal purposes, no matter whether the firm is controlled directly, indirectly, through foreign jurisdictions or multiple corporate vehicles (Knobel, Harari, and Meinzer 2018).

This paper, based on the research carried out within the EU co-funded project DATACROS², addresses these two gaps. On the one side, we carry out the first large-scale empirical analysis of the opacity of the ownership structure of more than 3 million firms, investigating the relationship between their ownership characteristics and the presence of enforcement/sanction measures on them or on their owners. On the other side, we shift the focus from the *who* controls the companies

² DATACROS (www.transcrime.it/datacros) is a research project co-funded by the European Commission, DG Home Affairs. It is coordinated by Transcrime-Università Cattolica del Sacro Cuore (Italy) and sees the partnership of the AFA, the French anti-corruption agency, the Spanish police (Cuerpo Nacional de Policia) and of the Italian network of investigative journalists IRPI. DATACROS will end in February 2021.

to the *how* control takes place, testing whether the ownership structure of the firm could be measured and whether the obtained measures are helpful to identify cases at high ML/TF risk and to predict potential sanctions, freezing orders and enforcement measures.

For doing so, we develop three indicators – the BOC, BOS and BOU – which summarise three facets of the problem of corporate ownership opacity: firms’ anomalous complexity (BOC), ownership links with ‘risky’ jurisdictions (BOS) and the unavailability of the information on beneficial owners (BOU) respectively. Then, we validate these indicators through a wide array of advanced statistical models against evidence of sanctions and enforcement on companies and their owners in 9 European countries.

Results confirm that opacity and anomalies in corporate ownership are more frequent in criminal cases. Namely, that the developed opacity indicators have a strong predictive capability in identifying firms (and their owners) that are then targeted by enforcement and sanctions. In particular, we are able to correctly detect 81% of positive cases of sanctions on companies, 70% of enforcements on companies, 88% of sanctions on BOs, and 62% of enforcements on BOs. Results also suggest that, among the analysed countries, there are differences in terms of predictive power and relative relevance of each one of these corporate opacity metrics.

Apart from the advancements in terms of empirical AML/CTF research, the proposed models also have strong policy implications. On the one side, they may help public authorities (AML/CTF supervisors, FIUs, law enforcement, policy-makers at large) to better monitor how ML/TF risks distribute across regions and sectors, and to identify areas or group of companies at higher risk on which to devote further investigation. On the other side, they may support AML/CTF obliged entities such as banks, professionals and gaming companies, in improving the performance of their customer due diligence (CDD) activities and of suspicious transaction reporting (STRs), boosting the predictive power of banks’ risk assessment models and minimizing the number of false positives.

2. Literature Review

2.1. The opacity of corporate ownership structure and ML/TF risk

As mentioned in the introduction, there is a general agreement among researchers and institutions that opaque corporate structures may hide and facilitate money laundering and other financial crimes (van Duyne and van Koningsveld 2017; FATF 2014). But there is still no common agreement on what corporate *opacity* means from the perspective of AML/CTF research. On the basis of a review of the literature in this field, we can identify three different facets of opacity: 1. Anomalous complexity of ownership structures; 2. Ownership links with high-risk jurisdictions; and 3. Ownership links with opaque corporate vehicles. We will discuss these three factors here below.

2.1.1. Anomalous complexity of ownership structures

The first facet of opacity relates to the complexity of the ownership chain of certain firms. In some cases, beneficial owners control their firms through lengthy and complex chains, made of cross-links and several layers of intermediate holding firms with a so-called ‘Chinese boxes’ scheme (i.e. an individual controlling a company, which in turn controls another company which in turn controls another company, and so on). Both the mentioned ‘StAR’ report (van der Does de Willebois et al. 2011) and the BOWNET report (Riccardi and Savona 2013) provide numerous cases of complex ownership chains employed in money laundering, corruption and financial crime cases. Also recent works on organized crime infiltration in the economy (Transcrime 2018) provide various examples of this modus operandi, which is considerably frequent for certain criminal offences, such as VAT or tax fraud that are extremely facilitated by these ‘Chinese boxes’ structures (Borselli 2011; Hangacova and Stremy 2018). Other studies have shown the employment of complex corporate structures in collusive or corruptive behaviors in public procurements (Fazekas, Tóth, and King 2013). The figure below represents a case of complex ownership structure employed in a case of mafia infiltration in the Italian legitimate economy, which allowed representatives of a *Cosa Nostra* family to obtain the monopoly of logistic and security services provided to a leading food retailer and to the local tribunal (Transcrime 2018).

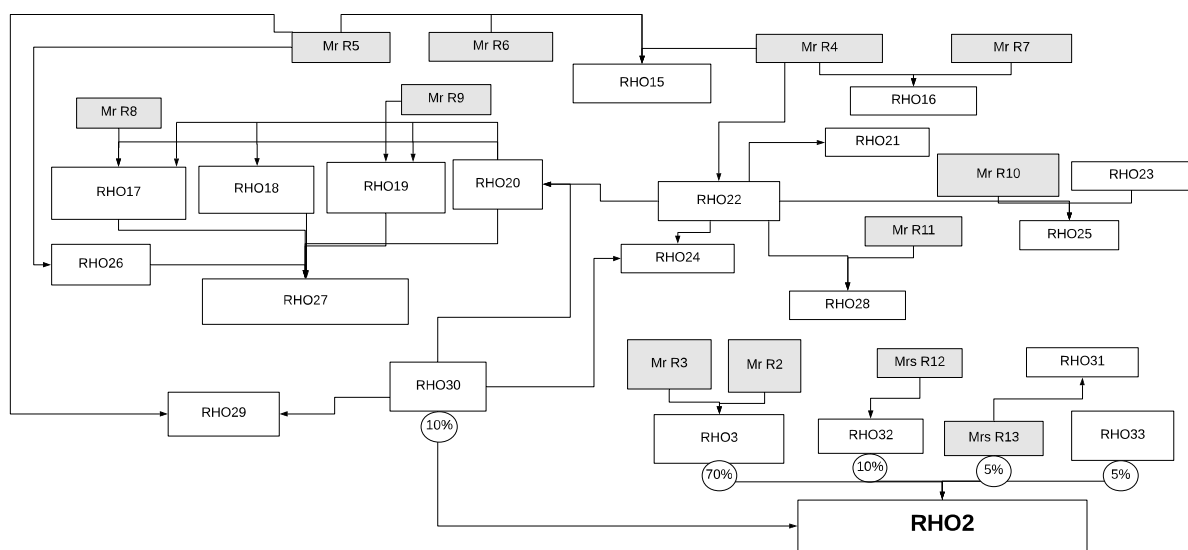


Figure 1: Ownership chain employed by Cosa Nostra affiliates to control company Rho2. All names of companies and owners involved in the case have been anonymized. Source: Transcrime, 2018.

However, it has to be stressed that complexity itself is not enough to label a company as ‘risky’ from the ML/TF perspective, because in some cases complex structures could be justified by the wide extent of the business activity carried out by the firm, its geographical reach or its financial strategy (e.g. controlling an asset or a controlled entity through debt leverage). For example, in some capital-intensive sectors such as the pharmaceutical industry, mining, and oil & gas, firms generally show a more complex structure characterized by the presence of various intermediate holding entities. This may be due also to tax optimization purposes (especially when holdings are registered in low tax jurisdictions), which may sometimes conceal more aggressive (and criminal) tax practices, but that are not necessarily related to money laundering.

Nonetheless, complex structures certainly pose a challenge to investigators willing to identify the natural persons ultimately lying behind a company, because they make it necessary to perform several queries on the company register – or on different business registers if the ownership chain spans cross-border. For this reason, some segments of the civil society, fostered by NGOs such as the Tax Justice Network and Transparency International, are pushing policy makers to introduce limits on the length of corporate ownership chains, or at least enhanced justifications in those cases in which individuals control firms through anomalous lengthy intermediate shareholdings (Tax Justice Network n.d.).

2.1.2. Ownership links with high-risk jurisdictions

Another facet of corporate ownership opacity relates to the existence of ownership links with high-risk jurisdictions – i.e. the presence of first level shareholders, intermediate shareholders or beneficial owners from countries with high levels of corporate and financial secrecy. The hypothesis is that if a beneficial owner controls a firm registered, say, in a country A with high level of transparency through another firm registered in a country B with low level of transparency, it would make it difficult for investigators to identify her/his identity (Aziani, Ferwerda, and Riccardi 2020).

The problem here is defining what ‘high-risk country’ means due to the lack of a general consensus around an agreed list of ‘secrecy’ jurisdictions. While the FATF and the EU carry out an in-depth analysis of beneficial ownership level of disclosure (through the evaluation of FATF Recommendations 24 and 25), several scholars have pointed out the inherent biases in FATF assessment methodologies (Halliday, Levi, and Reuter 2014; Levi, Reuter, and Halliday 2018; van Duyne and van Koningsveld 2017). In fact, evidence stemming from recent journalistic investigations, such as *Panama Papers* or *Paradise Papers*, points towards another direction, which is not correlated with FATF blacklist and greylist; while it is more aligned with the rankings of financial secrecy, such as those produced by the Tax Justice Network (Tax Justice Network 2018), or the blacklists and grey lists issued in the tax domain (e.g. the EU list of non-cooperative tax countries). However, these countries’ lists as well have their own methodological and political biases (see Riccardi 2020 for a review).

The debate on secrecy jurisdictions is never ending and falls outside the scope of this paper. What is worth stressing here is that criminals may make use of intermediate shareholding companies set up in secrecy jurisdictions in order to conceal their identity and the illicit source of their proceeds; and therefore that detecting and measuring these ownership links could help the identification of high risks of ML/TF.

2.1.3. Ownership links with opaque corporate vehicles

The third facet of corporate opacity is related to the employment, in the ownership chain, of legal entities (sometimes referred to as legal arrangements) which do not allow, by statute, the identification of beneficial owners, such as trusts, fiduciaries, foundations and some types of investment funds. The hypothesis is that criminals may control a firm through these legal entities in order to conceal their identity behind a legal veil (Riccardi and Savona 2013). Evidence about the employment of corporate vehicles in money laundering cases is wide, ranging from laundering of grand corruption to the recent Malaysia's 1MDB (Knobel 2019a; van der Does de Willebois et al. 2011).

The ML/TF vulnerabilities of legal arrangements have been frequently stressed by scholars and institutions, including the FATF (FATF 2006; 2010; OECD 2001; HM Revenue & Customs 2010). While the BO identification discipline has also evolved in several countries so as to find ways to identify, at least officially, beneficial owners also in the case of these legal arrangements (see, e.g., IADB and OECD 2019; HM Revenue & Customs 2010), these measures have been criticized by observers to be purely formal and not able to identify the 'real' BO. This is particularly true in the case of certain investment funds (Knobel 2019a) and foundation-types, such as the Dutch *stichting* (OECD 2019; Netherlands Chamber of Commerce - KVK 2020).

2.2. Research questions

The review of the literature - both the academic and institutional one – confirms the risks that corporate ownership opacity may pose in terms of ML/TF activities. But these alerts are backed by case-studies only, and are somehow the result of a global debate which, for the last 15 years, is constantly reiterating itself without providing further evidence in support. What do empirical data say, instead? Do they confirm the threat advocated by scholars and international AML/CTF authorities?

This paper addresses this issue. In particular, it aims at answering two research questions:

- (a) First, whether it is possible to *measure* corporate ownership opacity, and specifically whether it is possible to produce some indicators summarizing the three facets of opacity as emerging from the extant literature, that is anomalous complexity, ownership links with high-risk jurisdictions and ownership links with opaque legal arrangements;
- (b) Second, whether the newly developed indicators are positively associated with criminal activities, and are able to identify and predict firms (and their owners) then targeted by sanctions and enforcement measures.

We believe that responding these two questions is crucial both for advancing the empirical research in AML/CTF and for providing practitioners – both public authorities and obliged entities – with more effective tools to assess ML/TF risks and detect suspicious transactions and customers. Expanding the empirical evidence on which the global debate on corporate opacity is based could only improve the quality of the discussion and avoid unintended consequences – such as de-risking towards certain countries, industry categories or legal forms.

3. Methodology

3.1. Data

3.1.1. Data sources

For producing the indicators of ownership opacity at company level, data is collected from two sources (business ownership data and country blacklists data). A third source (sanctions and enforcement measures) is employed to validate the indicators. The three data sources are illustrated in detail here below.

Business ownership data

Information on 3,064,089 million companies in nine European countries (Belgium, Cyprus, Spain, France, United Kingdom, Italy, Luxembourg, Malta and Netherlands) is retrieved from *Orbis Europe*, a dataset provided by *Bureau van Dijk*. The exploited dataset provides a snapshot of European businesses during the month of June, 2019. In order to guarantee cross-country and cross-sector comparability, only limited companies with information on the ownership structure are considered in the analysis. We decided to focus only on Europe because of (a) the more homogeneous data coverage ensured by Bureau van Dijk Orbis in Europe, and the higher data quality, and (b) the more homogenous AML/CTF regulatory framework (all the mentioned jurisdictions fall below the EU 4th and 5th AML Directive, at least until the Brexit will substantially change the AML/CTF regime of the United Kingdom).

The extracted data include information on geographic location of the company (country, region, postcode), business sector of activity (NACE Rev.2 classification), business size (total assets, revenue, number of employees), and full ownership structure, i.e. all entities that connect the company with its *Beneficial Owners* (BOs), which are the individuals who ultimately own or control it. It is possible to identify them by reconstructing the ownership chain of a company, until finding natural persons with shareholding above a certain level. While most AML regulations worldwide define as beneficial owners the natural persons holding, directly or indirectly, at least 25% of the company share capital, or exerting a substantial control over the company in any case, for the purpose of this study we take a more prudential approach, setting a minimum threshold of 10% of shareholding at each level of the company ownership chain, value that is in line with the current AML debate and proposals suggested by various institutions and NGOs working in this field (see e.g. Knobel 2019; Riccardi and Milani 2018).

The coverage of company information across different countries is shown in *Figure 2*. As it can be seen, Italy, Spain and the United Kingdom are those countries characterised by the highest number of firms included in our dataset. While this data is aligned with the official number of registered businesses at the local company register, it may be also biased by the different degree of coverage ensured by Bureau van Dijk Orbis across countries³.

³ For more details on this, see Aziani, Ferwerda, and Riccardi 2020.

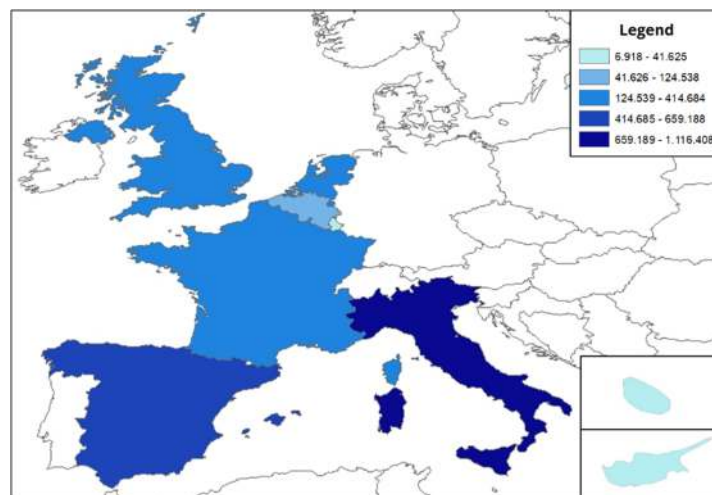


Figure 2: Geographic distribution of companies considered in the study

Country blacklists data

For operationalising the concept of ‘high-risk jurisdictions’ discussed in the previous section, we employ the following official black and grey lists. We decided to rely on official black and grey list in order to align with the official public statements and the current practice of banks and public authorities when assessing the risk of the jurisdiction where a counterpart or a client come from. However, we decided not to limit only to AML/CTF blacklists and greylists, but to expand it to the non-cooperative tax domain since, as discussed thoroughly in the previous section, the issue of corporate ownership opacity entails both AML and tax concerns. Moreover, being the focus of the analysis Europe, we decided to employ the EU list in the tax domain, while for AML/CFT we opted for the FATF one (as the newly European list of high-risk third countries does not cover, by definition, EU member states, neither potentially). In particular, we employed the following lists:

- *Tax domain*: the EU black and grey list of non-cooperative jurisdictions for tax purposes (updated 8th November, 2019), which groups countries that encourage abusive tax practices, which erode EU member states' corporate tax revenues (European Commission 2020);
- *AML/CTF domain*: Financial Action Task Force lists of non-cooperative jurisdictions in the global fight against money laundering and terrorist financing (October, 2019 statement). In particular, two lists are included:
 - *Call for action* (or so-called ‘Blacklist’) identifies countries that are evaluated by FATF as non-cooperative in the global fight against money laundering and terrorist financing, flagging them as "Non-Cooperative Countries or Territories" (NCCTs);
 - *Other monitored jurisdictions* (or so-called ‘Greylist’) identifies jurisdictions that have strategic AML/CFT deficiencies for which they have developed an action plan together with the FATF (for more information, see FATF 2017; 2019).

<i>List</i>	<i>Countries included</i>
EU blacklist of non cooperative jurisdictions for tax purposes (08/11/2019)	American Samoa, Fiji, Guam, Oman, Trinidad and Tobago, United States Virgin Islands, Vanuatu, Samoa
EU grey list of non cooperative jurisdictions for tax purposes (08/11/2019)	Anguilla, Antigua and Barbuda, Armenia, Australia, Bahamas, Barbados, Bermuda, Bosnia and Herzegovina, Botswana, Belize, British Virgin Islands, Cape Verde, Cayman Islands, Cook Islands, Curacao, Jordan, Maldives, Marshall Islands, Mongolia, Montenegro, Morocco, Namibia, Nauru, Niue, Palau, Saint Kitts and Nevis, Saint Lucia, Seychelles, Swaziland, Thailand, Turkey, Vietnam
FATF AML blacklist (October 2019 statement) - Call for action	Iran, Democratic People's Republic of Korea
FATF AML greylist (October 2019 statement) - Other monitored jurisdictions	Bahamas, Bouvet Island, Cambodia, Ghana, Iceland, Mongolia, Palau, Papua New Guinea, Tajikistan, Tunisia, Yemen, Zimbabwe

Table 1: Black and greylists of countries used for the study

Sanction and enforcement data

Information on companies and business owners included in one or more sanction lists or associated with enforcement cases is used to validate the indicators, as it will be illustrated below. This information is obtained from Lexis Nexis *WorldCompliance* for the nine countries at issue. In particular, we consider:

- *Enforcement lists*: companies and individuals with enforcement provisions (e.g. arrests, judgments) and court filings around the world, collated by Lexis Nexis from various sources including National law enforcement reports, press releases and other public authorities' statements. For the purpose of this paper, all categories of crimes and predicate offences covered by Lexis Nexis were taken into account;
- *Sanction lists*: companies and individuals included in one or more of the global screening and sanction lists issued by the following institutions: the European Union, the US Office of Foreign Assets Control (OFAC), the United Nations, the Bank of England, the US Federal Bureau of Investigation, the US Bureau of Industry and Security (BIS) and others⁴.

Data is obtained at company level and further matched to business ownership data. For the purposes of the study, we consider the following risk flags:

- *Company sanction*: dummy variable indicating if a company appears in *WorldCompliance* (as of June, 2019) as reported in one or more sanction lists;

⁴ For more information, see <https://risk.lexisnexis.com/global/en/products/worldcompliance-data>.

- *Company enforcement*: dummy variable indicating if a company appears in *WorldCompliance* as reported in one or more enforcement lists or court filings (e.g. seizure of the company or of its shares, administrative measure or ban on the companies, and others);
- *BOs sanction*: dummy variable indicating if a company has one or more beneficial owners appearing in *WorldCompliance* as reported in one or more sanction lists;
- *BOs enforcement*: dummy variable indicating if a company has one or more beneficial owners appearing in *WorldCompliance* as reported in one or more enforcement lists or court filings (e.g. arrests, final judgments, personal precautionary measures, freezing orders, others).

The geographical distribution of the above-mentioned variables is shown in *Figure 3*.

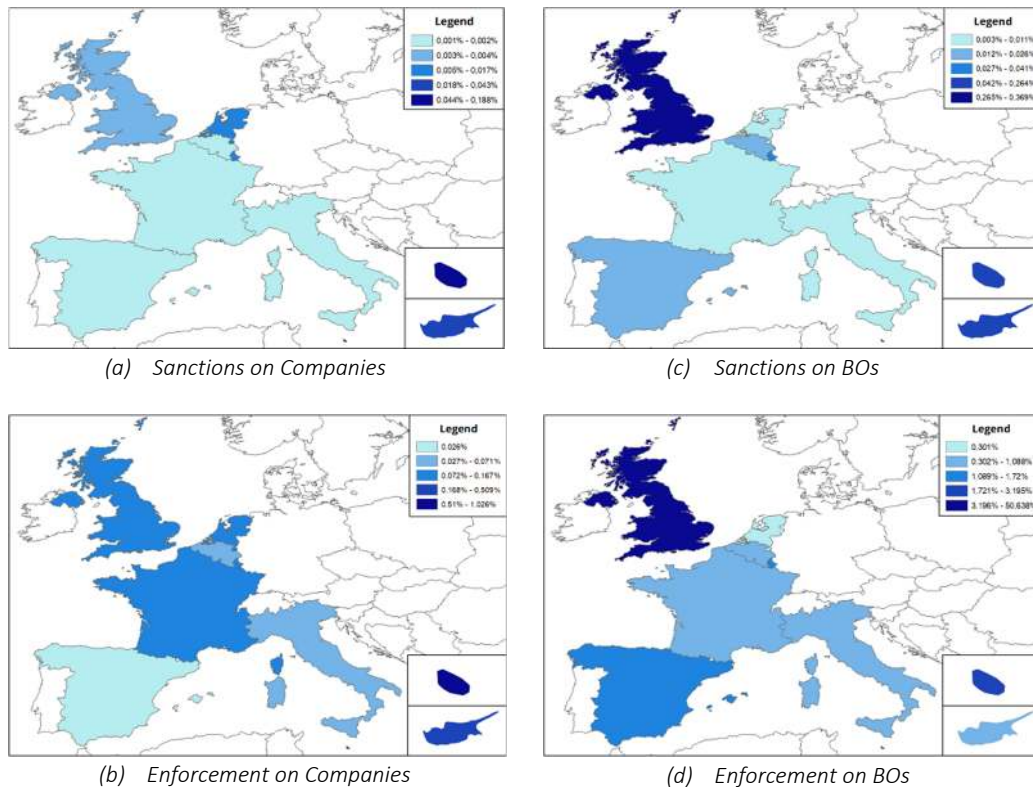


Figure 3: Geographical distribution of negative evidence data among the considered countries. From left to right, in the top row (a) data on sanctioned companies and (b) data on enforced companies, and in the bottom row (c) data on sanctioned BOs and (d) data on enforced BOs.

3.1.2. Ownership opacity risk indicators

For each company in the extracted dataset, we define and further calculate three risk indicators for investigating the three different aspects of ownership opacity which have been discussed in the previous section, including anomalous complexity of ownership, ownership links with blacklisted countries, and unavailability of information on beneficial owners.

For this end, we reconstruct the full ownership chain connecting each company to its business owners, even when deployed across borders. Accordingly, for each company we identify entities with shareholding above 10% at each level, until we reach the ultimate individual beneficiaries at the top of the chain (i.e. the BOs). When it is not possible to identify individuals at the top a chain, we name the last shareholder *Other Ultimate Beneficiary* (OUB). All entities separating a company from its BOs and/or OUBs are referred to as *intermediate shareholders* (INT).

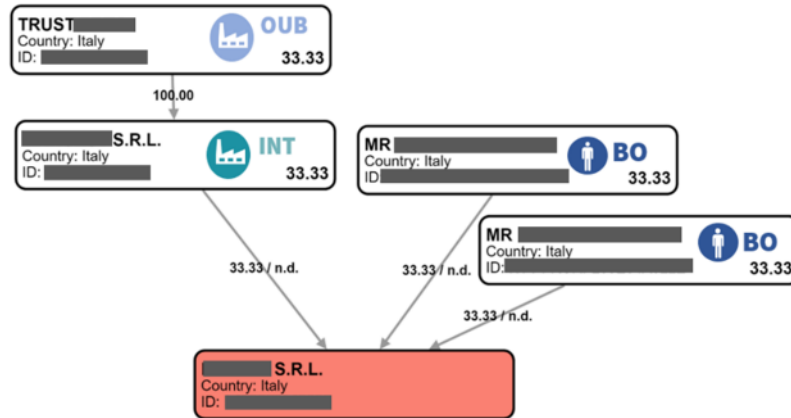


Figure 4: Illustration of different actors across ownership structures, including Beneficial Owners (BOs), Other Ultimate Beneficiaries (OUBs) and intermediate shareholders (INTs).

Beneficial ownership complexity (BOC)

The first indicator aims at measuring the risk related to the complexity of the company's ownership structure in order to facilitate the identification of possible anomalies, which is done by statistically comparing the characteristics of the company ownership chain with those of its peers operating in the same business sector and based on nearby geographical locations.

For the purposes of this study, we quantify complexity by looking at what is known as the *BO distance* that represents the number of steps separating a company from its ultimate BOs/OUBs. When BO distance is equals to 1, then the company is directly controlled by its BOs/OUBs, i.e. without intermediate shareholders or layers. The rationale for this approach is that the greater the BO distance, the higher the level of complexity of the company structure. Then, the more complex the chain is, the more difficult it is to trace beneficial owners, which in turn represents a greater challenge for investigative authorities and a higher risk that a company may be used to conceal criminal profits and/or individuals. Nevertheless, it is important to underline that the complexity of ownership is not anomalous *per se*. Anomalies arise when companies have an unnecessarily complex corporate structure with respect to the nature of their activity and their dimension. Therefore, we separate the sample into groups of peer companies, the so-called *Peer Groups*, which are groups of companies active in the same business sector and with a comparable dimension⁵. Next, within each Peer Group, we plot the distribution of BO distance and further segment it into 5 classes using a k-means hierarchical clustering algorithm. In this way, every company in the sample is assigned an anomalous complexity risk score (in a range of 1 to 5), hereinafter referred to as *BOC*, which indicates anomalous levels of complexity with respect to companies of the same Peer Group.

⁵ The segmentation in Peer Groups was taken from *Orbis Europe*. Peer Groups assemble together companies active in the same sector (NACE rev.2 classification) and belonging to the same dimensional class. Within each sector, five dimensional classes are defined based on the distribution of turnover and total assets: Very Small (VS), Small (SM), Medium (ME), Large (LA), Very Large (VL).

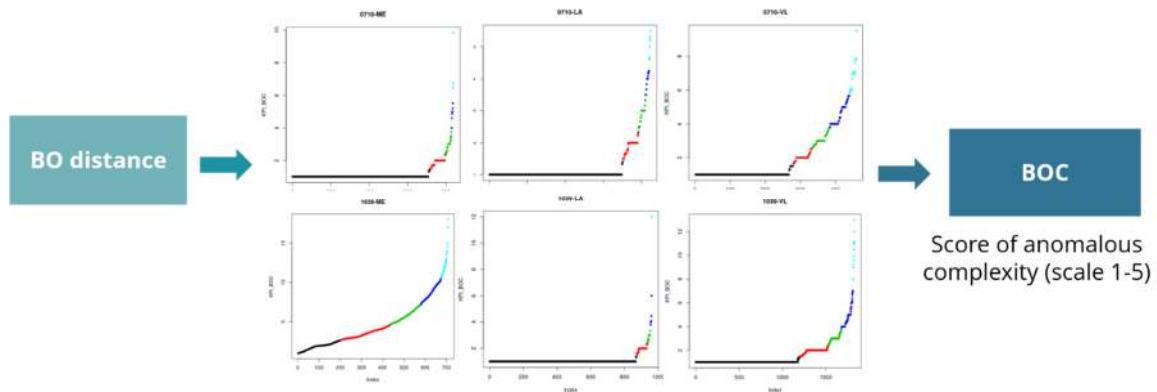


Figure 5: Clustering performed across the whole sample within each Orbis Peer Group defined by company sector and dimension.

Beneficial ownership secrecy (BOS)

The second indicator measures the risk related to the exposure of a corporate chain to entities located in high-risk jurisdictions, as defined above, i.e. as countries included in the FATF blacklist and greylists and in the EU list of non-cooperative tax jurisdictions. When a company is located in a country with such vulnerabilities, it is more difficult to carry out financial investigations and trace beneficial owners. As discussed in the previous section, there is a higher risk that these companies could be employed as intermediate shareholders in the ownership chain and to hide individuals and proceeds related to criminal activities.

In order to understand the level of business ownership links with secrecy jurisdictions, we match data on the nationality of the owner (if natural person) or location of the firm (if the owner is a legal person) with the above mentioned black and grey lists of risky jurisdictions using the whole sample. For each company, we count the number of BOs, OUBs and INTs linked to risky countries, and then divide by the total number of entities in the ownership chain. In doing so, we calculate for each company i the following proportion:

$$h_i = \frac{S_i^{BL}}{S_i}$$

where S_i = total number of BOs, OUBs or INTs in the ownership chain and S_i^{BL} = total number of BOs, OUBs or INTs in the ownership chain from a blacklisted country.

Same as BOC, we partition the metric into 5 classes by employing a k-means hierarchical clustering algorithm using the whole sample. In this way, each company is assigned a secrecy risk score (in a range of 1 to 5) called BOS, which quantifies the strength of links with secrecy jurisdictions.

Beneficial ownership unavailability (BOU)

The third and last indicator aims to quantify the risk related to the unavailability of information of company owners related to the employment of legal arrangements. As stressed by numerous institutions, researchers and NGOs, the more difficult it is to unequivocally identify the BOs, the greater the risk that a company could be used for concealing illicit activities (Knobel 2019b; Tax Justice Network 2018; Riccardi and Milani 2018).

In some cases, the identification of BOs of a company is not possible, which can happen for several reasons. For instance, if a company's shareholding is highly fragmented (such as for many listed

companies), there might be no individuals owning more than 10% of the shares. In other cases, legal arrangements can be used to hide the identity of individuals at the top of the ownership chain. The first is not a high-risk scenario, while the second, on the basis of the literature above conducted, may entail higher risks. Therefore, we want to distinguish the two situations. Consequently, for each company we identify, when possible, all OUBs that are trusts, fiduciaries, foundations or investment funds, hence do not allow the identification of BOs. In a like manner, we calculate for each company i the following proportion:

$$k_i = \frac{T_i^{TRUST}}{T_i}$$

where T_i = total number of OUBs in the ownership chain and T_i^{TRUST} = total number of OUBs in the ownership chain that are trusts, fiduciaries, foundations or investment funds.

Again, we divide the metric into 5 classes using a k-means hierarchical clustering algorithm. In this way, each company in the sample is assigned an unavailability risk score (in range of 1 to 5), hereinafter referred to as *BOU*, which will serve as indicator of the prevalence of trusts and/or other opaque corporate vehicles within the ownership chains.

The resulting clusters for the different indicators present heterogeneous distributions, which can be observed in *Figure 6*. Interesting patterns are further identified when analyzing clusters distribution by country and sector of activity (*Table 2* and *Table 3* -). It is worth noting that for all considered indicators, lower cluster values indicate lower risk hence more attention should be paid to entities categorized with values 5 and 4.

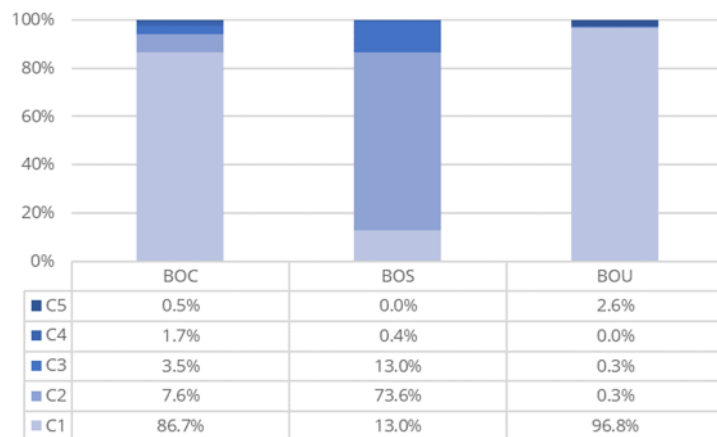


Figure 6: Distribution of clusters of the different indicators BOC, BOS and BOU

In the case of complexity, it can be seen that while most observations are classified as low risk, there is a remarkable presence of companies categorized as ‘high-risk’ from Netherlands, Luxembourg, Malta and Belgium. Also a dominant presence of risky businesses operating in sectors D (*Electricity, gas, steam and air conditioning supply*) and K (*Financial and insurance activities*), and in a lesser extent in sectors B (*Mining and quarrying*), H (*Transporting and storage*) and R (*Arts, entertainment and recreation*). In addition, a notable proportion of companies from Cyprus are categorized as high risk with regard to the secrecy indicator. As for unavailability of BO information, we observe an overall 2.6% of observations classified as very risky; businesses operating in sector K (*Financial and insurance activities*) stand out for the high presence of this risk factor (15.2%), as well as companies

operating in the Netherlands (23.5%), probably for the extensive use of Dutch *stichting*, foundation-type legal arrangements often also used for controlling for-profit limited or unlimited companies (OECD 2019; Netherlands Chamber of Commerce - KVK 2020).

Country	BOC					BOS					BOU				
	C1	C2	C3	C4	C5	C1	C2	C3	C4	C5	C1	C2	C3	C4	C5
BE - Belgium	82.4%	8.0%	5.1%	3.4%	1.1%	1.0%	88.5%	9.0%	1.5%	0.0%	97.4%	0.9%	0.8%	0.0%	0.9%
CY - Cyprus	81.1%	12.3%	4.0%	2.1%	0.5%	0.4%	3.1%	87.8%	8.4%	0.3%	96.1%	0.6%	1.1%	0.1%	2.1%
ES - Spain	89.5%	6.4%	2.6%	1.3%	0.3%	0.5%	97.8%	1.7%	0.1%	0.0%	99.2%	0.1%	0.1%	0.0%	0.6%
FR - France	77.5%	12.3%	6.2%	3.0%	0.9%	0.8%	95.7%	3.4%	0.1%	0.0%	98.6%	0.1%	0.4%	0.0%	0.8%
GB - United Kingdom	92.6%	4.3%	1.7%	1.0%	0.4%	93.1%	5.2%	1.6%	0.1%	0.0%	99.4%	0.3%	0.2%	0.0%	0.0%
IT - Italy	90.7%	6.3%	2.1%	0.8%	0.1%	0.3%	97.4%	2.1%	0.2%	0.0%	99.3%	0.1%	0.2%	0.0%	0.4%
LU - Luxembourg	77.1%	11.9%	5.8%	3.8%	1.4%	1.3%	15.0%	81.4%	2.3%	0.1%	91.2%	2.2%	2.0%	0.2%	4.3%
MT - Malta	62.2%	16.5%	12.8%	7.1%	1.4%	1.0%	6.7%	89.6%	2.5%	0.1%	95.9%	1.3%	1.0%	0.0%	1.8%
NL - Netherlands	72.0%	13.2%	8.1%	4.8%	1.9%	0.3%	2.1%	97.0%	0.6%	0.0%	74.4%	1.0%	1.0%	0.1%	23.5%

Table 2 : Distribution of clusters of the different indicators by country

NACE (Rev.2 classification)	BOC					BOS					BOU				
	C1	C2	C3	C4	C5	C1	C2	C3	C4	C5	C1	C2	C3	C4	C5
A - Agriculture, forestry and fishing	85.3%	6.7%	4.9%	2.5%	0.6%	7.2%	81.1%	11.6%	0.1%	0.0%	97.2%	0.1%	0.2%	0.0%	2.5%
B - Mining and quarrying	72.7%	14.5%	8.2%	3.7%	1.0%	16.3%	71.0%	11.8%	0.8%	0.0%	95.2%	0.8%	2.3%	0.0%	1.8%
C - Manufacturing	84.4%	8.6%	4.6%	1.9%	0.5%	7.3%	84.1%	8.4%	0.2%	0.0%	98.0%	0.2%	0.4%	0.0%	1.4%
D - Electricity, gas, steam and air conditioning supply	58.4%	20.4%	11.6%	7.6%	1.9%	5.8%	82.2%	11.5%	0.5%	0.0%	96.4%	0.3%	0.5%	0.0%	2.8%
E - Water supply	76.2%	15.3%	5.4%	2.4%	0.7%	12.9%	77.2%	9.8%	0.1%	0.0%	97.0%	0.3%	0.4%	0.0%	2.2%
F - Construction	90.7%	5.5%	2.0%	1.5%	0.3%	11.6%	81.4%	6.8%	0.3%	0.0%	98.7%	0.1%	0.2%	0.0%	1.1%
G - Wholesale and retail trade	88.7%	6.8%	2.8%	1.3%	0.5%	6.9%	80.1%	12.5%	0.5%	0.0%	97.9%	0.2%	0.3%	0.0%	1.6%
H - Transporting and storage	83.1%	9.7%	4.0%	2.3%	0.9%	10.8%	75.1%	13.4%	0.7%	0.0%	97.2%	0.3%	0.3%	0.0%	2.2%
I - Accommodation and food service activities	91.0%	6.2%	2.1%	0.6%	0.1%	10.5%	82.1%	7.0%	0.4%	0.0%	98.9%	0.1%	0.1%	0.0%	0.9%
J - Information and communication	86.5%	7.9%	3.5%	1.8%	0.4%	25.1%	63.5%	11.2%	0.2%	0.0%	98.0%	0.3%	0.4%	0.0%	1.3%
K - Financial and insurance activities	75.4%	9.7%	8.4%	5.0%	1.5%	8.0%	46.6%	44.2%	1.1%	0.1%	82.9%	1.0%	1.0%	0.0%	15.0%
L - Real estate activities	86.8%	9.0%	3.0%	0.9%	0.3%	13.2%	77.5%	9.1%	0.2%	0.0%	97.1%	0.2%	0.3%	0.0%	2.4%
M - Professional, scientific and technical activities	87.2%	8.1%	2.8%	1.5%	0.4%	18.8%	62.5%	18.5%	0.3%	0.0%	96.2%	0.3%	0.4%	0.0%	3.2%
N - Administrative and support service activities	85.5%	8.9%	3.4%	1.7%	0.5%	24.5%	59.5%	15.5%	0.5%	0.0%	96.4%	0.5%	0.6%	0.1%	2.4%
O - Public administration and defence	93.4%	3.8%	1.1%	1.5%	0.1%	63.9%	27.5%	8.0%	0.3%	0.3%	99.0%	0.3%	0.4%	0.0%	0.4%
P - Education	91.2%	6.6%	1.5%	0.6%	0.1%	19.7%	68.9%	11.3%	0.1%	0.0%	97.6%	0.1%	0.3%	0.0%	2.0%
Q - Human health and social work activities	86.8%	6.9%	3.8%	1.8%	0.8%	22.1%	65.0%	12.8%	0.1%	0.0%	96.8%	0.2%	0.3%	0.0%	2.7%
R - Arts, entertainment and recreation	86.1%	6.7%	4.4%	1.9%	0.9%	17.1%	72.4%	10.3%	0.2%	0.0%	97.7%	0.1%	0.2%	0.0%	2.0%
S - Other services activities	92.5%	4.5%	1.7%	0.8%	0.5%	32.6%	59.6%	7.5%	0.4%	0.0%	98.9%	0.1%	0.2%	0.0%	0.8%
T - Activities of households as employers	94.6%	1.9%	2.4%	0.8%	0.3%	98.4%	1.3%	0.3%	0.0%	0.0%	99.0%	0.7%	0.2%	0.0%	0.0%
U - Activities of extraterritorial organisations and bodies	88.0%	6.0%	3.9%	1.4%	0.7%	10.1%	70.9%	17.4%	1.5%	0.0%	99.0%	0.0%	0.3%	0.0%	0.7%

Table 3 - Distribution of clusters of the different indicators by business sector defined by NACE Rev.2 classification

3.1.3. Sampling

The main concern with the collected and processed dataset relates to the imbalance of all the target variables under study. In all considered scenarios, sanctions and enforcement are rare events leading to a very large proportion of non-criminal instances in the data. The ratio of the small to the large class is quite drastic in all these scenarios especially considering that less than 0.1% of entities reports negative evidence for the first three dependent variables, that is *Sanction on companies*, *Enforcement on companies* and *Sanction on BOs*, and only 7.8% in the case of *Enforcement on BOs*.

It is worth noting that when breaking down the analysis by country, it can be observed (Table 4) that higher proportions of offences on companies relate to businesses in Malta and Cyprus, and to a lesser extent to businesses in Luxembourg and Netherlands. For offences on BOs, on the contrary, we observe that the most predominant country is the United Kingdom, with an exceptional 51% of companies having at least one BO targeted by enforcement measures. This surprising result could be explained not only by the higher percentage of ‘criminals’ (or presumed so) in British firms, but also by the higher quality of coverage by Lexis Nexis in the UK, which could bring to a relatively higher presence of negative events on UK firms. The asymmetry between the number of enforced

BOs from the United Kingdom and the rest of the countries generates complications for the models as inferences are constructed based on the most common class (i.e., companies from the United Kingdom) disregarding any additional information from companies located in other locations. As result, we obtain useless models incapable of distinguishing risky companies from less predominant countries. For this reason, we conduct a separate analysis for *Enforcement on BOs*, first considering all countries excluding GB, and second, considering only companies from the United Kingdom.

Country	Company sanction	Company enforcement	BOs sanction	BOs enforcement
All countries	0.004%	0.087%	0.064%	7.801%
BE - Belgium	0.002%	0.071%	0.026%	0.768%
CY - Cyprus	0.043%	0.509%	0.264%	0.913%
ES - Spain	0.001%	0.026%	0.021%	1.522%
FR - France	0.001%	0.123%	0.011%	0.990%
GB - United Kingdom	0.004%	0.167%	0.369%	50.638%
IT - Italy	0.001%	0.050%	0.005%	1.088%
LU - Luxembourg	0.017%	0.119%	0.041%	1.720%
MT - Malta	0.188%	1.026%	0.173%	3.195%
NL - Netherlands	0.014%	0.129%	0.003%	0.301%

Table 4: Percentage of WorldCompliance flags on the total number of companies – first as a whole, then by country. Volumes by country are calculated considering the total number of companies within each state.

As for the present study, imbalanced data represents a relevant issue for modelling and prediction as statistical learning models are highly affected by the non-criminal class. The discrimination between the two classes becomes a very difficult task leading to techniques ignoring criminal offences and opting for a strategy of classifying all instances as licit. Consequently, overall accuracy measures and true negative rates show an outstanding predictive performance that only reflect the underlying uneven class distribution (Chawla, Japkowicz, and Kołcz 2004; He and Garcia 2009). Nevertheless, the methods are totally ineffective in detecting events of money laundering and related offences. In order to circumvent the imbalance problem, we first perform a selection of a more proportionate sample in terms of licit and illicit cases, and then we evaluate results of predictive models using more suitable measurements of performance. The proposed strategy will help to enhance the discriminatory power of the statistical techniques as well as the interpretation of relevant results.

Accordingly, the sampling strategy employed in the present study is based on the undersampling of the majority class (i.e. non-sanctioned/non-enforced observations, hereinafter referred to as *negatives*) so to randomly match the number of observations with the minority class (i.e. sanctioned/enforced observations, hereinafter referred to as *positives*). The random pairing exercise, which is summarized in *Table 5*, is performed in all scenarios except in the case of Enforcement on BOs for the United Kingdom. As mentioned before, the number of criminal cases from the UK compared to the rest of the countries is extremely large, so we do not perform undersampling in this case but take the entire sample as it is already balanced in its original form. In the case of Sanctions on companies, on the contrary, we are forced to perform an oversampling of positive cases since the number of offences is exceptionally low, with only 117 observations. Consequently, we first select all criminal observations and split them by country, and then we perform a random sample with replacement on each stratum. The replacement option allows us to

extract samples of any volume, thus giving us the possibility to artificially increase the size of the positive datasets, in this case by 10.

	Company sanction	Company enforcement	BOs sanction	BOs enforcement GB excl.	BOs enforcement GB only
Total num obs	3,064,089	3,064,089	3,064,089	2,649,405	414,684
Negative events	117	2,677	1,950	29,056	209,987
Oversampling rate	x10	N/A	N/A	N/A	N/A
Oversampled obs	1,170	N/A	N/A	N/A	N/A
Undersampled obs	1,170	2,677	1,950	29,056	204,697
Sample size	2,340	5,354	3,900	58,112	414,684

Table 5: Information on sampling and final samples for the different targets

Moreover, and as it can be observed in Table 6 and Table 7, all samples obtained from the random undersampling and oversampling exercises show very similar country and sector specific distribution as their original datasets, which is consequence of large sample sizes. In light of this information, we can say that all resulting samples are, in fact, representative and that can be further used for generalization of results.

	Company sanction				Company enforcement		BOs sanction		BOs enforcement GB excl.	
	Negatives (undersampling)		Positives (oversampling)		Negatives (undersampling)		Negatives (undersampling)		Negatives (undersampling)	
	Dataset	Sample	Dataset	Sample	Dataset	Sample	Dataset	Sample	Dataset	Sample
BE - Belgium	4.1%	4.4%	1.7%	2.1%	4.1%	4.7%	4.1%	4.4%	4.7%	4.8%
CY - Cyprus	1.4%	1.8%	15.4%	15.1%	1.4%	1.0%	1.4%	1.2%	1.6%	1.5%
ES - Spain	21.5%	22.5%	6.8%	5.9%	21.5%	22.2%	21.5%	22.3%	24.8%	24.8%
FR - France	12.9%	11.1%	4.3%	5.3%	12.9%	12.8%	12.9%	12.1%	14.9%	14.9%
GB - United Kingdom	13.5%	14.1%	14.5%	16.1%	13.5%	13.3%	13.5%	14.1%	-	-
IT - Italy	36.4%	35.0%	9.4%	8.6%	36.4%	35.7%	36.5%	36.3%	42.1%	42.4%
LU - Luxembourg	1.1%	1.6%	5.1%	5.3%	1.1%	1.2%	1.1%	1.1%	1.3%	1.3%
MT - Malta	0.2%	0.0%	11.1%	11.4%	0.2%	0.3%	0.2%	0.2%	0.3%	0.3%
NL - Netherlands	8.8%	9.5%	31.6%	30.3%	8.8%	8.7%	8.9%	8.5%	10.3%	10.2%

Table 6: Distribution of observations of original datasets and the resulting samples, for all target variables and by country.

	Company sanction				Company enforcement		BOs sanction		BOs enforcement GB excl.	
	Negatives (undersampling)		Positives (oversampling)		Negatives (undersampling)		Negatives (undersampling)		Negatives (undersampling)	
	Dataset	Sample	Dataset	Sample	Dataset	Sample	Dataset	Sample	Dataset	Sample
A - Agriculture, forestry and fishing	1.7%	1.7%	1.7%	0.9%	1.7%	1.2%	1.7%	1.6%	1.8%	1.8%
B - Mining and quarrying	0.2%	0.1%	0.2%	4.5%	0.2%	0.3%	0.2%	0.5%	0.2%	0.2%
C - Manufacturing	10.6%	10.2%	10.6%	7.8%	10.6%	10.5%	10.6%	10.6%	11.3%	11.4%
D - Electricity, gas, steam and air conditioning supply	1.1%	0.7%	1.1%	0.0%	1.1%	0.9%	1.1%	1.1%	1.1%	1.1%
E - Water supply	0.4%	0.3%	0.4%	0.0%	0.4%	0.4%	0.4%	0.5%	0.4%	0.5%
F - Construction	14.3%	12.9%	14.3%	4.4%	14.3%	13.7%	14.3%	13.3%	14.5%	14.4%
G - Wholesale and retail trade	21.0%	22.6%	21.0%	14.1%	21.0%	21.9%	21.0%	19.7%	22.5%	22.6%
H - Transporting and storage	3.3%	3.0%	3.3%	13.2%	3.3%	3.1%	3.3%	3.8%	3.3%	3.3%
I - Accommodation and food service activities	5.9%	5.7%	5.9%	0.0%	5.9%	6.2%	5.9%	6.1%	6.0%	6.1%
J - Information and communication	4.9%	4.1%	4.9%	2.3%	4.9%	5.0%	4.9%	5.0%	4.2%	4.2%
K - Financial and insurance activities	5.9%	7.2%	5.9%	32.6%	5.9%	6.1%	5.9%	6.3%	6.3%	6.1%
L - Real estate activities	8.6%	7.9%	8.6%	5.5%	8.6%	8.6%	8.6%	9.8%	8.6%	8.6%
M - Professional, scientific and technical activities	10.1%	10.8%	10.1%	4.9%	10.1%	10.2%	10.1%	10.3%	9.5%	9.4%
N - Administrative and support service activities	5.4%	5.6%	5.4%	9.0%	5.4%	5.2%	5.4%	5.3%	4.7%	4.7%
O - Public administration and defence	0.0%	0.1%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
P - Education	0.9%	0.9%	0.9%	0.3%	0.9%	0.6%	0.9%	0.9%	0.8%	0.8%
Q - Human health and social work activities	1.8%	2.2%	1.8%	0.6%	1.8%	2.0%	1.8%	1.9%	1.6%	1.6%
R - Arts, entertainment and recreation	1.7%	1.9%	1.7%	0.0%	1.7%	1.6%	1.7%	1.9%	1.6%	1.6%
S - Other services activities	1.8%	1.7%	1.8%	0.0%	1.8%	2.1%	1.8%	1.5%	1.3%	1.4%
T - Activities of households as employers	0.4%	0.5%	0.4%	0.0%	0.4%	0.4%	0.4%	0.3%	0.0%	0.0%
U - Activities of extraterritorial organisations and bodies	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%

Table 7: Distribution of observations of original datasets and the resulting samples, for all target variables and by economic sector.

3.2. Models

For the next task of modelling the phenomena at issue, we opt for simple models in terms of number of predictors as we are most interested in validating risk indicators and identify differences across locations and sectors of activity. Accordingly, we consider negative evidence on sanctions and enforcement as dependent variables, and further use the three proposed indicators as predictors. Country of location and economic sector of activities are employed as controls. An illustration of variables used in the methods in the figure below.



Figure 7: Variables included in the modelling. We consider 3 predictors (BOC, BOS, BOU), 2 controls (country and economic sector) and 4 target variable (sanctions on companies, enforcement on companies, sanctions on BOs and enforcement on BOs)

3.2.1. Model assessment

As mentioned in the previous section, the dataset we have to cope with to investigate the phenomenon of interest is unbalanced with respect to the target variables of sanctions and enforcements on companies and BOs. The issue with this kind of data is that resulting models are excellent at identifying licit entities, but useless on flagging corrupted ones. In the context of money laundering, flagging a company as non-risky when, in fact, it is (false negative), represents a much higher misclassification cost for society and governments than flagging a company as risky when it is not (false positive). Therefore, overall accuracy and true negative rate for imbalance datasets are not sufficient for assessing models' performance. In line with this, we choose to partially renounce to correctly classify licit companies for the sake of correctly categorize illicit ones. Consequently, other metrics are also taken into consideration including true positive rate, AUC and precision.

It is worth mentioning that the calculation of all these measurements are based on the confusion matrix shown below.

	Predicted positives	Predicted negatives	Total
Real positives	TP	FN	P
Real negatives	FP	TN	N

Table 8: Confusion matrix

Model assessment metrics are described next along with the formulas used to calculate them when appropriate.

Measures of accuracy

- *True positive rate (TPR)*: also called *sensitivity*. It assesses the capacity to classify criminal cases correctly. It is then calculated as the proportion of true positive cases compared to the actual number of criminal observations.

$$TPR = \frac{TP}{P}$$

- *True negative rate (TNR)*: also called *specificity*. It evaluates the ability to classify good entities correctly. As such, it is computed as the proportion of true negative compared to the number of licit observations.

$$TNR = \frac{TN}{N}$$

- *Overall accuracy*: it measures the ability to differentiate between licit and illicit entities correctly. It is calculated as the proportion of true positive and true negative cases compared to the total number of observations.

$$\text{Overall Accuracy} = \frac{TP + TN}{P + N}$$

- *Area under the curve (AUC)*: it summarizes the quality of the Receiver Operating Characteristic (ROC) curve as a single value. The ROC curve plots the true positive rate against the false positive rate for different threshold settings, thus illustrating the false positive rate that must be accepted to obtain a particular level of true positive rate.
- *Precision*: it measures the proportion of true positive cases compared to the total number of observations predicted as illicit, that is, the proportion of detections that are actually positive.

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

It is worth noting that all the aforementioned accuracy measures are probabilities, hence will always be positive numbers ranging from zero to one; so the closer to the unit, the better is the model. For the purpose of assessment then, we will tend to favor models with relatively high values of most of these measurements as they evidence good predictive performance.

The hold-out set approach

The learning process of a statistical model requires the use of data in order to estimate the associated parameters and to evaluate performance. With these two purposes in mind, we adopt the *hold-out* set approach. This technique is a simple but very powerful technique based on the random split of the dataset at issue into two sets that are equally balanced in terms of the target variable. The first set, representing the 80% percent of all the data, will be attributed to what is called the *training set*, and the remainder 20% to what is known as the *hold-out set*.

What is sought during the training phase is to minimize the costs of misclassifications by fitting the algorithm to the data at hand. As result, we obtain optimal parameters that best fit the data used for this purpose hence performing extremely well when classifying. Nonetheless, the true predictive power of the models can only be assessed using a different set of data that has not been used during the training set and therefore has not influenced the estimates. Accordingly, to properly evaluate the performance of the models we use this last set to calculate all accuracy measures defined in the previous section.

3.2.2. Machine learning methods

The discrimination and further identification of sanctioned/enforced and non-sanctioned/non-enforced entities in all considered scenarios, correspond to what is known in statistical jargon as a binary classification problem. Thus, machine learning models estimate first the probability that an observation is positive and then, if this probability is higher than 50%, classify the instance as illicit. An exhaustive analysis on the probability threshold has resulted in very small changes to the classifications and subsequent performance of the models, hence it is decided to maintain the traditional 50%-approach as to avoid confusion and ease the understanding of the outcomes.

Accordingly, this study assesses the performance of several machine learning models in the identification of money laundering offences related to both companies and BOs. First, commonly used techniques, including Logistic Regression, Naïve Bayes Classifier and a stacking alternative based on both, are employed as benchmark framework. Then, we implement more advanced but still comprehensible tree-based algorithms, including decision trees, bagged trees and random forests. All models have been optimized and fitted for the next step of classification, and further tested using *ad hoc* samples and accuracy metrics properly described in the previous section of model assessment.

In addition, it is worth mentioning that the implementation of most models involves the estimation of optimal hyperparameters governing such methods. In the case of logistic regression, for instance, different $L-1$ and $L-2$ regularized forms have also been considered but not further adopted as they never outperformed the non-regularized alternative (James et al. 2014). As for the tree-based methods, 5-fold cross-validation is used for the tuning of parameters such as tree depth, minimum sample size of end-nodes and maximum number of random predictors, among others.

All machine learning models are employed as implemented in the Scikit-Learn library (Pedregosa et al. 2011).

3.2.3. Robustness Analysis

We further analyze how results change when cases from a certain country or business sector are excluded. This robustness check should give us some evidence on two main affairs. First, it can be used to investigate whether our indicators are significant in different criminal scenarios, and second, it can be employed to quantify the relative relevance of the indicators in different countries and sectors. For this end, we use the whole dataset, that is, no splitting of training and testing sets, and only the indicators as predictors, as we are controlling for countries and sectors manually. We then fit a logistic regression using the information of all countries and sectors, and set it as a benchmark for comparing the remaining logistic models that are fitted excluding one country/sector at a time. By doing this, we validate the significance of the coefficients related to the indicators, both in general terms as well as when breaking down the analysis by group of interest. Furthermore, we can calculate the relative difference between the odds-ratio when considering all observations and when removing a specific group, and further use these values to assess the relative importance of the indicators in different countries and sectors.

4. Results

4.1. Models performance

Results on the performance of different machine learning methods can be observed in *Table 9*. As mentioned before, all accuracy measures are calculated using the *hold-out* set approach described in the previous section.

<i>Models</i>	<i>TPR</i>	<i>TNR</i>	<i>Overall Accuracy</i>	<i>AUC</i>	<i>Precision</i>
<i>Company sanction</i>					
Logistic Regression	0.833	0.872	0.853	0.931	0.867
Naive Bayes Classifier	0.739	0.850	0.795	0.875	0.832
Stacking	0.833	0.872	0.853	0.931	0.867
Decision Tree	0.876	0.846	0.861	0.919	0.851
Bagged Trees	0.910	0.778	0.844	0.918	0.804
Random Forest	0.752	0.868	0.810	0.922	0.850
<i>Company enforcement</i>					
Logistic Regression	0.679	0.729	0.704	0.785	0.715
Naive Bayes Classifier	0.621	0.721	0.671	0.751	0.691
Stacking	0.687	0.729	0.708	0.784	0.717
Decision Tree	0.769	0.634	0.701	0.731	0.678
Bagged Trees	0.763	0.634	0.698	0.759	0.676
Random Forest	0.729	0.662	0.696	0.766	0.684
<i>BOs sanction</i>					
Logistic Regression	0.879	0.851	0.865	0.896	0.855
Naive Bayes Classifier	0.872	0.854	0.863	0.880	0.856
Stacking	0.879	0.851	0.865	0.895	0.855
Decision Tree	0.869	0.856	0.863	0.879	0.858
Bagged Trees	0.856	0.841	0.849	0.890	0.843
Random Forest	0.851	0.859	0.855	0.885	0.858
<i>BOs enforcement GB excl.</i>					
Logistic Regression	0.615	0.564	0.589	0.634	0.585
Naive Bayes Classifier	0.566	0.568	0.567	0.597	0.567
Stacking	0.615	0.564	0.589	0.634	0.585
Decision Tree	0.520	0.675	0.598	0.639	0.615
Bagged Trees	0.515	0.670	0.592	0.640	0.609
Random Forest	0.578	0.610	0.594	0.649	0.597
<i>BOs enforcement GB only</i>					
Logistic Regression	0.548	0.522	0.535	0.550	0.540
Naive Bayes Classifier	0.510	0.545	0.527	0.537	0.535
Stacking	0.548	0.522	0.535	0.550	0.540
Decision Tree	0.874	0.164	0.524	0.533	0.518
Bagged Trees	0.874	0.164	0.524	0.533	0.518
Random Forest	0.550	0.523	0.537	0.554	0.542

Table 9: Prediction accuracy of different models for the different target variables

It can be seen that results are dissimilar across different models and target variables. In terms of response, best performance of the proposed methods is achieved when predicting sanctions, both of BOs and of companies. On the contrary, lowest performance is obtained when predicting cases of enforced BOs, either excluding the United Kingdom or considering it in isolation, which suggests that this phenomenon in particular is relatively difficult to detect and forecast compared to the other cases.

As for the machine learning models, it is worth noting that logistic regression is always amongst the best alternatives. More advanced tree-based techniques show similar predictive capability compared to benchmark models, except in the case of enforcements on BOs in the United Kingdom where decision tree and bagged trees outperform by almost 33%. Moreover, it can be said that the stacking model performs virtually the same as logistic regression, which suggests that the Naïve Bayes classifier does not stand out in any of the situations studied.

Lastly, true positive rates of logistic regression models indicate very promising performance as we are able to correctly detect 83% of positive cases of sanctions on companies, 68% of enforcements on companies, 88% of sanctions on BOs, and 62% of enforcements on BOs⁶.

4.2. Indicators performance

Understanding the relevance and the predictive ability of the three proposed indicators is crucial for validating their usefulness as red flags of money laundering. We certainly prefer machine learning models that perform well, but we are much more interested in quantifying the importance of different aspects of opacity, that is complexity, secrecy and unavailability of BO information, in determining the occurrence of criminal offences so they can be ultimately used for prevention and detection.

Information on the performance of the indicators can be retrieved from several of the statistical methods we have implemented in this study. In particular, we can assess logistic regression outcomes in the form of significance of coefficients and of odds-ratio, as well as random forests by analyzing predictors contribution by means of variable importance.

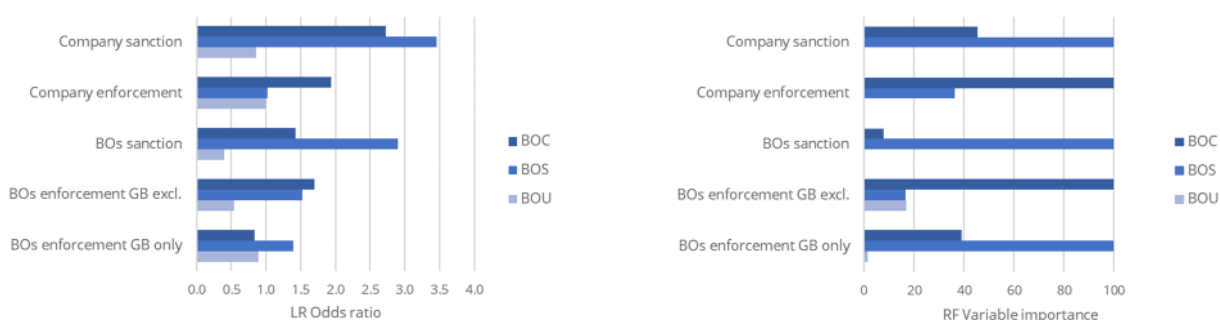


Figure 8: Indicators performance on the basis of logistic regression odds ratio (left) and of random forests variable importance (right)

In regard to statistical significance, we find that most coefficients of logistic regression models are significant at a 0.001 level except for the second case of enforcement on companies where neither BOU nor BOS are significant. Furthermore, it is observed in both graphs from Figure 8, that BOC is

⁶ Excluding the United Kingdom (GB).

notably important for detecting most offences, in particular enforcement cases. As for the BOS indicator, we find evidence on its ability to predict sanction offences, both on companies and BOs, as well as enforcement on BOs in the United Kingdom. Lastly, BOU seems to be the least relevant indicator in all considered cases.

4.3. Indicators assessment

As a result of the robustness analysis proposed in the previous section, we are able to obtain a more detailed overview of the relevance of indicators by country and by economic sector. It is worth noting that values presented in both tables below are ratios illustrating the relative effect, in terms of value change of logistic regression coefficients, obtained by excluding observations from specific segments. Therefore, positive ratios indicate countries or sectors where the relative contribution of a specific indicator is greater; these cases are highlighted in red-toned colors in *Table 10* and *Table 11*, facilitating the identification of potential red flags⁷.

Regarding the analysis by country, some interesting trends can be observed by looking at the indicators. In the first place, we can see (*Table 10*) that most indicators seem robust and consistent across countries. Nevertheless, some countries with major and systematic upward/downward variation of the results can be identified. For instance, we observe that corporate ownership complexity (BOC) is particularly relevant in Italy for identifying enforcement and sanction cases involving companies, while in Spain it is a particularly strong factor for the identification of sanctions and enforcements against BOs. These results could be related to the distribution of complex companies in these two countries: in Italy and Spain, companies are relatively more often directly controlled by their BOs⁸, and are less likely to have maximum score of anomalous complexity (0.1% and 0.3% of total population, see *Table 2* for details). Therefore, few cases of anomalous complexity among companies with negative evidence can drive up the results. In terms of secrecy, we observe that the BOS indicator plays a very relevant role for both enforcements on companies and sanctions on BOs in the United Kingdom. It is also interesting to notice that in the Netherlands, the unavailability of BO information appears to be a predominant factor for most negative evidence related to BOs. This is a very relevant result, considering that 23% of limited Dutch companies are ultimately controlled by a trust, a fiduciary or a fund that does not allow for the identification of a BO (see *Table 2*). This is most likely connected to the extended domestic use of Dutch foundations, schemes also known as *stichting*: legal arrangements that are similar to foundations (in the sense that they are designed for managing philanthropical and non-profit objectives), but in the Netherlands they are also heavily used to control for-profit limited or unlimited firms. In fact, for *stichting*, it is not mandatory to disclose the 'owners' of the foundations (Netherlands Chamber of Commerce - KVK 2020), which, in turn, makes them suitable for hiding the identity of ultimate beneficiaries of controlled firms. Not by chance, the vulnerability of *stichting* to money laundering, tax evasion and criminal activities has been stressed by various agencies, such as OECD that in its 2019 report, underlined that: "*Foundations in the Netherlands are not systematically required to keep identity information concerning all beneficiaries. An obligation should be established in both*

⁷ More detail on the result of robustness analysis by country in Appendix 1 and by economic sector in Appendix 2.

⁸ The average BO distance in Italy (1.15) and Spain (1.21) is below EU average (1.23), meaning that domestic companies in these two countries show - on average - lower levels of complexity than in other countries.

the European Netherlands and the Caribbean Netherlands for foundations to keep identity information concerning all beneficiaries” (OECD 2019).

Country	Company sanction			Company enforcement			BOs sanction			BOs enforcement		
	BOC	BOS	BOU	BOC	BOS	BOU	BOC	BOS	BOU	BOC	BOS	BOU
BE - Belgium	0.00	0.02	-0.12	0.00	0.10	1.07	0.01	-0.02	0.03	-0.04	-0.59	0.02
CY - Cyprus	0.03	0.13	0.36	-0.09	-3.71	-2.28	0.58	-0.36	0.15	0.00	0.08	-0.01
ES - Spain	0.02	0.13	-0.33	0.03	0.17	1.68	0.45	0.10	-0.03	0.16	0.59	-0.06
FR - France	-0.16	0.10	-0.04	-0.01	-0.41	-0.26	-0.10	0.02	-0.13	-0.19	-0.35	-0.05
GB - United Kingdom	0.02	-2.09	-1.96	0.04	12.58	2.37	-0.55	2.50	-0.87	-	-	-
IT - Italy	0.18	0.08	-0.72	0.14	0.30	1.84	0.38	0.19	-0.18	-0.03	-1.82	0.04
LU - Luxembourg	-0.01	0.04	-0.03	-0.03	-0.05	0.33	0.00	-0.03	0.08	0.01	-0.22	-0.01
MT - Malta	-0.11	0.34	1.41	0.01	-0.71	-0.66	-0.04	-0.03	0.01	-0.01	-0.19	0.00
NL - Netherlands	0.12	0.00	-1.68	-0.06	-0.81	-5.33	0.03	0.08	0.33	0.05	3.19	0.30

Table 10: Robustness analysis by country

In like manner, we can see interesting patterns emerging from the sector-specific analysis summarized in Table 11. While most results hold stable across sectors, we see that complexity is a major determinant of enforcement and sanction on BOs in sector K (*Financial and insurance activities*), one of the sector with the greatest presence of companies with maximum score of anomalous complexity (1.5%, see Table 3). We also see that in sector K the secrecy indicator BOS plays an important role for identifying sanctions on companies and enforcements on BOs⁹. Lastly, we also observe that unavailability of BO information (BOU) is an important factor for the identification of most types of negative evidence in sector G (*Wholesale and retail trade*), and of sanctions on companies operating in sector H (*Transporting and storage*). It is worth mentioning that to properly address the interpretation of some of these trends, the analysis would require a much deeper sector-level, country-level and country-sector-level specific exploration that goes beyond the scope of this study.

NACE (Rev.2 classification)	Company sanction			Company enforcement			BOs sanction			BOs enforcement GB excl.			BOs enforcement GB only		
	BOC	BOS	BOU	BOC	BOS	BOU	BOC	BOS	BOU	BOC	BOS	BOU	BOC	BOS	BOU
A - Agriculture, forestry and fishing	0.00	-0.01	0.29	0.00	0.01	-0.10	-0.05	0.01	-0.01	-0.01	0.02	0.00	0.00	0.00	0.02
B - Mining and quarrying	0.02	0.05	0.02	-0.01	-0.09	0.15	0.00	0.00	0.01	0.01	0.00	0.00	0.01	-0.01	0.02
C - Manufacturing	0.03	-0.08	0.13	0.04	-0.03	0.08	0.06	0.00	-0.01	0.00	-0.04	0.00	-0.04	-0.01	-0.08
D - Electricity, gas, steam and air conditioning supply	-0.02	0.01	-0.05	-0.02	-0.13	-0.22	-0.08	0.00	-0.02	0.03	0.07	0.00	-0.02	0.00	0.01
E - Water supply	0.00	0.00	-0.01	0.00	0.14	0.12	-0.02	0.00	0.00	0.00	0.02	0.00	0.01	0.00	0.00
F - Construction	0.04	-0.08	-0.41	0.00	1.19	0.78	0.01	0.05	-0.11	-0.02	0.05	-0.03	0.00	0.09	-0.07
G - Wholesale and retail trade	0.04	0.20	0.78	0.08	0.08	1.33	-0.24	-0.05	0.02	-0.12	-0.46	0.03	0.05	-0.01	0.20
H - Transporting and storage	-0.06	0.20	1.15	0.01	-0.26	-0.34	-0.12	-0.04	-0.01	-0.02	-0.02	0.01	0.04	0.04	0.01
I - Accommodation and food service activities	-0.02	0.03	-0.08	0.01	0.19	0.46	0.01	0.05	0.05	-0.05	-0.28	0.00	0.02	0.01	0.13
J - Information and communication	0.00	0.02	0.04	0.01	-0.02	-0.01	0.21	-0.01	0.00	0.02	0.05	0.00	0.02	0.00	-0.09
K - Financial and insurance activities	-0.07	0.29	-12.71	-0.04	-0.86	-1.90	0.47	-0.15	0.24	0.06	0.27	0.01	-0.07	-0.04	0.05
L - Real estate activities	0.01	-0.26	0.21	-0.02	0.34	0.21	-0.01	0.06	-0.08	0.01	0.08	-0.01	0.10	-0.07	0.02
M - Professional, scientific and technical activities	-0.02	0.01	0.54	-0.02	-0.68	-0.53	-0.04	-0.01	-0.07	0.05	0.28	0.01	-0.12	0.02	-0.08
N - Administrative and support service activities	0.02	-0.40	0.14	-0.01	0.07	0.13	-0.09	-0.02	0.11	0.03	-0.11	0.01	0.00	-0.02	0.06
O - Public administration and defence	0.00	0.00	0.00	0.00	0.00	0.14	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
P - Education	0.00	0.01	-0.02	0.00	0.05	0.04	-0.01	0.01	0.00	0.00	0.00	0.00	-0.01	0.01	-0.03
Q - Human health and social work activities	0.01	-0.04	-0.12	-0.01	0.04	0.06	0.00	0.03	0.00	-0.01	0.04	-0.01	-0.01	-0.02	-0.10
R - Arts, entertainment and recreation	0.00	0.02	-0.06	0.00	-0.21	-0.37	0.06	0.02	0.00	0.00	-0.01	0.00	0.01	0.01	-0.04
S - Other services activities	0.00	0.00	-0.08	0.01	-0.03	-0.03	-0.06	0.01	-0.10	-0.01	-0.04	0.00	0.01	0.00	-0.02
T - Activities of households as employers	-0.01	0.02	0.00	0.00	-0.08	-0.02	0.01	0.01	0.00	0.00	0.00	0.00	0.00	0.00	-0.01
U - Activities of extraterritorial organisations and bodies	-	-	-	0.00	-0.01	-0.01	-	-	-	0.00	0.00	0.00	0.00	0.00	0.00

Table 11: Robustness analysis by economic sector (NACE Rev.2 classification)

⁹ Excluding the United Kingdom (GB).

5. Conclusions

The opacity of corporate ownership has become a central issue in the global AML/CTF domain for the last 15 years. Numerous measures have been implemented worldwide in order to increase the transparency of firms and their owners, first of all the establishment of BO registers. However, despite all the emphasis around the issue, the empirical evidence in support of this concern is poor, and mainly limited to some case-studies, while large-scale analyses were missing. The present paper, resulting from the research activity of the EU-funded Project DATACROS, coordinated by Transcrime-Università Cattolica del Sacro Cuore, has addressed this gap. It has proposed an innovative analytical approach for (a) measuring the opacity of corporate ownership through a set of risk indicators and (b) improving the risk assessment of money laundering and the detection of 'bad entities' (i.e. firms at high risk of being involved in criminal activity and being targeted by sanctions and enforcement) through these indicators.

This was done through an empirical analysis of 3,064,089 limited companies registered in nine European countries, including Belgium, Cyprus, Spain, France, United Kingdom, Italy, Luxembourg, Malta and Netherlands. Based on business ownership information and country blacklist data, we constructed three risk indicators related to (i) anomalous complexity of ownership *BOC*, (ii) ownership links with blacklisted countries *BOS*, and (iii) unavailability of information on beneficial owners *BOU*. In addition, we obtained target variables from sanction and enforcement lists.

Because of the imbalanced problem arising from the rare events of sanctions and enforcements, we performed a sampling strategy based on the undersampling of non-criminal offences and oversampling of criminal cases when little evidence of offences is available. We then used these samples to train and test different machine learning models for all the target variables at issue, and to further assess predictive performance not only of models but also of risk indicators.

Different results were obtained from the proposed methodology. As for the predictive capability of machine learning models, we see that best performance is achieved when predicting sanctions, which suggests enforcement cases are more difficult to detect and forecast. This may be explained by the wide variety of enforcement measures covered by the employed data source (Lexis Nexis World Compliance), which encompasses arrests, final judgments, freezing orders, confiscation orders, administrative measures and bans; and it may be also biased by the diverse capacity of the employed source to monitor and collect enforcement cases in different countries, while for sanctions the screening is easier since sanction lists are often made publicly available by the issuing institution.

In addition, we find logistic regression models to be always amongst the best alternatives and that the related estimated coefficients are significant in almost all scenarios. Moreover, their associated true positive rates indicate very promising predictive performance as we are able to correctly detect 83% of positive cases of sanctions on companies, 68% of enforcements on companies, 88% of sanctions on BOs, and 62% of enforcements on BOs¹⁰.

Interesting patterns emerge when breaking down the analysis by country and by economic sector. On the one hand, we find that results are robust across countries and sectors. On the other hand,

¹⁰ Excluding the United Kingdom (GB).

we find that some risk indicators show a relatively higher impact in some specific contexts. For instance, complexity plays a major role in countries such as Italy and Spain, probably for the distribution of complex companies in the two countries. The unavailability of BO information appears to be particularly relevant for identifying negative in the Netherlands, likely because of the extensive use of Dutch foundations (*stichting*) for controlling local companies. As for secrecy, we find it plays an important role in the United Kingdom, as well as for companies operating in sectors K (*Financial and insurance activities*) and F (*Construction*).

These results have both strong research and policy implications. In terms of research, they provide empirical support to the wide literature highlighting the relationship between corporate opacity and money laundering – and in general criminal activities. First, it goes beyond the extant literature which is based on a limited number of case studies. This paper, for the first time, conducts a systematic analysis of several millions firms demonstrating that this relationship holds also in ‘large amounts’: firms with a more complex or opaque ownership structure (and their owners) are more likely to be involved in criminal activities, i.e. to be targeted by sanctions or enforcement measures.

Second, it shifts the attention from the *who*, i.e. the natural persons controlling an entity, to the *how*, i.e. the way individuals govern firms. It provides an innovative way to operationalize the concept of ownership opacity and anomalies, through three metrics which could be easily adapted and applied to millions of firms worldwide.

In terms of policy implications, the models proposed by this paper provide users with newly tools to improve the risk assessment and the identification of high-risk firms. Both public authorities and private actors could benefit from the employment of such tools. First, AML/CTF supervisory agencies, FIUs and law enforcement agencies could use these models to better screening the companies under their jurisdictions in order to filter out those deserving in-depth investigation and audit. This activity, which is exactly what advocated in terms of *risk assessment*, may help these authorities to better allocate their resources, staff, instruments where risks are higher.

However, also private entities could benefit from the employment of the risk indicators here proposed – and of the related predictive models. In particular, AML/CTF obliged entities such as banks, financial institutions, professionals and gaming companies could make use of these indicators to improve the profiling of their clients and the effectiveness of the customer due diligence (CDD) activity carried out within their daily AML/CTF operations. Usually, AML/CTF obliged entities are suggested by relevant regulators and supervisory agencies a set of anomaly indicators to be taken into account when conducting CDD: but these are usually quite generic and still to be operationalized in actual terms. The indicators proposed here are ready-to-use and easily adaptable to most regulated entities’ datasets of clients (at least, most banks’ databases). The significant predictive power of these models may help the capacity of obliged entities to detect actual high-risk cases, reduce the number of false positives in suspicious transaction reports (STRs) and improve cost-efficiency of their activity, especially when dealing with large number of clients to be screened.

Appendix 1: Robustness analysis by country

		<i>Company sanction</i>									
		All	BE out	CY out	ES out	FR out	GB out	IT out	LU out	MT out	NL out
BOC		2.655*** (0.067)	2.666*** (0.069)	2.589*** (0.07)	2.593*** (0.072)	3.112*** (0.077)	2.594*** (0.078)	2.236*** (0.071)	2.67*** (0.068)	2.959*** (0.068)	2.372*** (0.08)
BOS		2.094*** (0.07)	2.071*** (0.07)	1.909*** (0.071)	1.906*** (0.069)	1.947*** (0.07)	9.821*** (0.123)	1.976*** (0.072)	2.036*** (0.071)	1.631*** (0.071)	2.09*** (0.076)
BOU		0.944 (0.074)	0.938 (0.074)	0.964 (0.075)	0.927 (0.075)	0.942 (0.076)	0.844** (0.078)	0.906 (0.074)	0.943 (0.075)	1.024 (0.074)	0.858 (0.154)
		<i>Company enforcement</i>									
		All	BE out	CY out	ES out	FR out	GB out	IT out	LU out	MT out	NL out
BOC		2.093*** (0.039)	2.096*** (0.041)	2.243*** (0.041)	2.052*** (0.042)	2.115*** (0.043)	2.033*** (0.042)	1.884*** (0.043)	2.136*** (0.04)	2.08*** (0.039)	2.195*** (0.045)
BOS		0.929 (0.046)	0.936 (0.047)	0.707*** (0.051)	0.941 (0.046)	0.901** (0.047)	2.344*** (0.072)	0.95 (0.047)	0.925* (0.047)	0.882*** (0.047)	0.875** (0.053)
BOU		1.026 (0.042)	0.998 (0.043)	1.089** (0.043)	0.982 (0.043)	1.033 (0.044)	0.965 (0.043)	0.978 (0.043)	1.017 (0.043)	1.044 (0.042)	1.179** (0.07)
		<i>BOs sanction</i>									
		All	BE out	CY out	ES out	FR out	GB out	IT out	LU out	MT out	NL out
BOC		1.361*** (0.059)	1.355*** (0.062)	1.139* (0.069)	1.184*** (0.064)	1.406*** (0.064)	1.615*** (0.07)	1.209*** (0.063)	1.362*** (0.06)	1.377*** (0.06)	1.348*** (0.062)
BOS		0.349*** (0.054)	0.343*** (0.055)	0.239*** (0.062)	0.389*** (0.054)	0.356*** (0.054)	4.881*** (0.115)	0.427*** (0.054)	0.339*** (0.055)	0.337*** (0.055)	0.378*** (0.056)
BOU		0.502*** (0.192)	0.511*** (0.188)	0.556*** (0.211)	0.492*** (0.187)	0.459*** (0.219)	0.275*** (0.376)	0.442*** (0.191)	0.529*** (0.187)	0.505*** (0.193)	0.63** (0.218)
		<i>BOs enforcement</i>									
		All	BE out	CY out	ES out	FR out	GB out	IT out	LU out	MT out	NL out
BOC		1.546*** (0.012)	1.573*** (0.012)	1.548*** (0.012)	1.441*** (0.013)	1.678*** (0.014)		1.567*** (0.014)	1.538*** (0.012)	1.555*** (0.012)	1.511*** (0.013)
BOS		0.824*** (0.023)	0.735*** (0.025)	0.837*** (0.025)	0.923*** (0.024)	0.77*** (0.025)		0.58*** (0.028)	0.79*** (0.025)	0.794*** (0.024)	1.528*** (0.029)
BOU		0.463*** (0.035)	0.469*** (0.036)	0.458*** (0.036)	0.441*** (0.041)	0.445*** (0.038)		0.475*** (0.038)	0.46*** (0.037)	0.462*** (0.036)	0.582*** (0.046)

Table 12: Robustness analysis: excluding sanction/enforcement cases from one country at a time

Appendix 2: Robustness analysis by economic sector

																						<i>Company sanction</i>									
	All	A out	B out	C out	D out	E out	F out	G out	H out	I out	J out	K out	L out	M out	N out	O out	P out	Q out	R out	S out	T out	U out									
BOC	2.655***	2.667***	2.596***	2.585***	2.711***	2.656***	2.544***	2.556***	2.803***	2.697***	2.651***	2.844***	2.638***	2.703***	2.613***	2.654***	2.65***	2.63***	2.662***	2.66***	2.684***	-									
	(0.067)	(0.068)	(0.067)	(0.072)	(0.068)	(0.067)	(0.069)	(0.073)	(0.068)	(0.069)	(0.069)	(0.085)	(0.069)	(0.071)	(0.069)	(0.067)	(0.067)	(0.067)	(0.068)	(0.068)	(0.068)	(0.068)	-								
BOS	2.094***	2.103***	2.024***	2.227***	2.076***	2.09***	2.226***	1.801***	1.807***	2.043***	2.07***	1.688***	2.529***	2.083***	2.81***	2.093***	2.082***	2.16***	2.068***	2.088***	2.067***	-									
	(0.07)	(0.07)	(0.069)	(0.076)	(0.07)	(0.07)	(0.073)	(0.071)	(0.072)	(0.07)	(0.071)	(0.077)	(0.077)	(0.072)	(0.08)	(0.07)	(0.07)	(0.071)	(0.07)	(0.07)	(0.07)	(0.07)	-								
BOU	0.944	0.96	0.945	0.951	0.941	0.944	0.922	0.987	1.009	0.94	0.947	0.456***	0.956	0.974	0.952	0.944	0.943	0.938	0.941	0.94	0.944	-									
	(0.074)	(0.074)	(0.074)	(0.075)	(0.074)	(0.074)	(0.074)	(0.077)	(0.074)	(0.074)	(0.074)	(0.073)	(0.197)	(0.077)	(0.078)	(0.075)	(0.074)	(0.073)	(0.074)	(0.074)	(0.074)	(0.074)	-								
																						<i>Company enforcement</i>									
	All	A out	B out	C out	D out	E out	F out	G out	H out	I out	J out	K out	L out	M out	N out	O out	P out	Q out	R out	S out	T out	U out									
BOC	2.093***	2.093***	2.112***	2.038***	2.125***	2.1***	2.086***	1.971***	2.084***	2.077***	2.08***	2.161***	2.128***	2.132***	2.114***	2.093***	2.086***	2.104***	2.087***	2.085***	2.094***	2.092***									
	(0.039)	(0.039)	(0.04)	(0.042)	(0.04)	(0.039)	(0.042)	(0.042)	(0.04)	(0.04)	(0.04)	(0.043)	(0.041)	(0.041)	(0.04)	(0.039)	(0.039)	(0.04)	(0.04)	(0.039)	(0.039)	(0.039)	(0.039)								
BOS	0.929	0.93	0.923*	0.927	0.92*	0.939	1.014	0.935	0.911**	0.942	0.927	0.872***	0.953	0.884**	0.934	0.929	0.932	0.932	0.914*	0.927	0.924*	0.928									
	(0.046)	(0.046)	(0.046)	(0.048)	(0.046)	(0.046)	(0.048)	(0.05)	(0.047)	(0.047)	(0.047)	(0.053)	(0.047)	(0.049)	(0.048)	(0.046)	(0.046)	(0.046)	(0.046)	(0.046)	(0.046)	(0.046)	(0.046)								
BOU	1.026	1.029	1.022	1.024	1.032	1.023	1.006	0.991	1.036	1.014	1.027	1.078	1.021	1.04	1.023	1.023	1.025	1.025	1.036	1.027	1.027	1.027									
	(0.042)	(0.042)	(0.042)	(0.044)	(0.042)	(0.042)	(0.043)	(0.045)	(0.042)	(0.042)	(0.042)	(0.049)	(0.043)	(0.044)	(0.043)	(0.042)	(0.042)	(0.042)	(0.042)	(0.042)	(0.042)	(0.042)	(0.042)								
																						<i>BOs sanction</i>									
	All	A out	B out	C out	D out	E out	F out	G out	H out	I out	J out	K out	L out	M out	N out	O out	P out	Q out	R out	S out	T out	U out									
BOC	1.361***	1.381***	1.363***	1.337***	1.395***	1.371***	1.358***	1.467***	1.413***	1.355***	1.276***	1.176**	1.365***	1.376***	1.398***	1.361***	1.366***	1.36***	1.336***	1.388***	1.359***	-									
	(0.059)	(0.06)	(0.059)	(0.062)	(0.06)	(0.059)	(0.061)	(0.066)	(0.061)	(0.06)	(0.062)	(0.067)	(0.06)	(0.062)	(0.061)	(0.059)	(0.059)	(0.059)	(0.059)	(0.059)	(0.059)	(0.059)	-								
BOS	0.349***	0.351***	0.348***	0.349***	0.35***	0.347***	0.368***	0.332***	0.336***	0.367***	0.346***	0.298***	0.373***	0.345***	0.342***	0.349***	0.351***	0.359***	0.355***	0.353***	0.352***	-									
	(0.054)	(0.054)	(0.054)	(0.056)	(0.054)	(0.054)	(0.056)	(0.06)	(0.056)	(0.055)	(0.056)	(0.059)	(0.055)	(0.057)	(0.058)	(0.054)	(0.054)	(0.054)	(0.054)	(0.055)	(0.054)	(0.054)	-								
BOU	0.502***	0.498***	0.505***	0.499***	0.495***	0.5**	0.464***	0.509***	0.499***	0.518***	0.501***	0.591**	0.473***	0.477***	0.543***	0.502***	0.502***	0.502***	0.502***	0.469***	0.502***	-									
	(0.192)	(0.193)	(0.191)	(0.192)	(0.193)	(0.192)	(0.212)	(0.202)	(0.195)	(0.184)	(0.207)	(0.206)	(0.208)	(0.218)	(0.189)	(0.192)	(0.192)	(0.191)	(0.191)	(0.213)	(0.191)	(0.191)	-								
																						<i>BOs enforcement GB excl.</i>									
	All	A out	B out	C out	D out	E out	F out	G out	H out	I out	J out	K out	L out	M out	N out	O out	P out	Q out	R out	S out	T out	U out									
BOC	1.546***	1.553***	1.543***	1.547***	1.523***	1.544***	1.563***	1.628***	1.557***	1.579***	1.532***	1.508***	1.537***	1.512***	1.529***	1.546***	1.546***	1.552***	1.549***	1.551***	1.546***	1.546***									
	(0.012)	(0.012)	(0.012)	(0.013)	(0.012)	(0.012)	(0.013)	(0.013)	(0.012)	(0.012)	(0.012)	(0.013)	(0.012)	(0.012)	(0.012)	(0.012)	(0.012)	(0.012)	(0.012)	(0.012)	(0.012)	(0.012)	(0.012)								
BOS	0.824***	0.827***	0.824***	0.817***	0.835***	0.827***	0.832***	0.754***	0.821***	0.781***	0.832***	0.868***	0.836***	0.87***	0.807***	0.824***	0.825***	0.83***	0.823***	0.818***	0.824***	0.824***									
	(0.023)	(0.024)	(0.023)	(0.024)	(0.023)	(0.023)	(0.025)	(0.028)	(0.024)	(0.024)	(0.024)	(0.025)	(0.024)	(0.025)	(0.024)	(0.023)	(0.023)	(0.024)	(0.024)	(0.024)	(0.023)	(0.023)	(0.023)								
BOU	0.463***	0.462***	0.463***	0.463***	0.463***	0.462***	0.453***	0.475***	0.465***	0.464***	0.462***	0.465***	0.457***	0.466***	0.465***	0.463***	0.462***	0.459***	0.463***	0.462***	0.463***	0.463***									
	(0.035)	(0.035)	(0.035)	(0.037)	(0.035)	(0.035)	(0.037)	(0.038)	(0.035)	(0.036)	(0.036)	(0.041)	(0.037)	(0.036)	(0.036)	(0.035)	(0.035)	(0.036)	(0.035)	(0.035)	(0.035)	(0.035)	(0.035)								
																						<i>BOs enforcement GB only</i>									
	All	A out	B out	C out	D out	E out	F out	G out	H out	I out	J out	K out	L out	M out	N out	O out	P out	Q out	R out	S out	T out	U out									
BOC	0.837***	0.837***	0.838***	0.832***	0.835***	0.838***	0.837***	0.845***	0.844***	0.841***	0.84***	0.827***	0.852***	0.82***	0.838***	0.838***	0.836***	0.837***	0.838***	0.839***	0.838***	0.838***									
	(0.006)	(0.006)	(0.006)	(0.007)	(0.006)	(0.006)	(0.007)	(0.007)	(0.006)	(0.006)	(0.006)	(0.007)	(0.007)	(0.007)	(0.007)	(0.006)	(0.006)	(0.007)	(0.006)	(0.006)	(0.006)	(0.006)	(0.006)								
BOS	1.418***	1.416***	1.423***	1.425***	1.419***	1.419***	1.372***	1.422***	1.397***	1.415***	1.418***	1.436***	1.455***	1.409***	1.426***	1.418***	1.414***	1.426***	1.415***	1.419***	1.42***	1.418***									
	(0.009)	(0.009)	(0.01)	(0.01)	(0.01)	(0.009)	(0.01)	(0.01)	(0.01)	(0.01)	(0.01)	(0.01)	(0.01)	(0.01)	(0.01)	(0.009)	(0.01)	(0.01)	(0.01)	(0.01)	(0.01)	(0.01)	(0.009)								
BOU	0.878***	0.88***	0.88***	0.87***	0.879***	0.87***	0.87***	0.901***	0.879***	0.893***	0.868***	0.884***	0.881***	0.87***	0.885***	0.878***	0.875***	0.867***	0.873***	0.876***	0.877***	0.878***									
	(0.022)	(0.022)	(0.022)	(0.023)	(0.022)	(0.022)	(0.023)	(0.023)	(0.022)	(0.023)	(0.023)	(0.024)	(0.024)	(0.024)	(0.023)	(0.022)	(0.022)	(0.022)	(0.022)	(0.022)	(0.022)	(0.022)	(0.022)								

Table 13: Robustness analysis: excluding sanction/enforcement cases from one economic sector at a time

References

- Aziani, Alberto, Joras Ferwerda, and Michele Riccardi. 2020. "Who Are Our Owners? Exploring the Cross-Border Ownership Links of European Businesses to Assess the Risk of Illicit Financial Flows." *European Journal of Criminology - in Course of Publication*, no. 1: 43.
- Borselli, F. 2011. "Organised VAT Fraud: Features, Magnitude, Policy Perspective." 106. *Questioni Di Economia e Finanza*. Roma: Banca D'Italia. http://www.bancaditalia.it/pubblicazioni/econo/temidi/td12/td868_12/en_td868/en_tema_868.pdf.
- Chawla, Nitesh, Nathalie Japkowicz, and Aleksander Kotcz. 2004. "Editorial: Special Issue on Learning from Imbalanced Data Sets." *SIGKDD Explorations*, June 1, 2004.
- Does de Willebois, Emile van der, Emiliy Halter M., Robert. A. Harrison, Ji Won Park, and J.C. Sharman. 2011. *The Puppet Masters: How the Corrupt Use Legal Structures to Hide Stolen Assets and What to Do About It*. The World Bank. <http://elibrary.worldbank.org/doi/book/10.1596/978-0-8213-8894-5>.
- Duyne, Petrus C. van, and T. J. van Koningsveld. 2017. "The Offshore World: Nebolous Finance." In *The Many Faces of Crime for Profit and Ways of Tackling It*.
- European Commission. 2020. "Taxation: EU List of Non-Cooperative Jurisdictions." 2020. <https://www.consilium.europa.eu/en/policies/eu-list-of-non-cooperative-jurisdictions/>.
- European Parliament, and Council of the European Union. 2015. *Directive (EU) 2015/849 of the European Parliament and of the Council of 20 May 2015. Official Journal of the European Union*. <http://eur-lex.europa.eu/legal-content/EN/TXT/?uri=celex%3A32015L0849>.
- FATF. 2006. "The Misuse of Corporate Vehicles, Including Trust and Company Service Providers." Paris: The Financial Action Task Force. <http://www.fatf-gafi.org/media/fatf/documents/reports/Misuse%20of%20Corporate%20Vehicles%20includ%20Trusts%20and%20Company%20Services%20Providers.pdf>.
- . 2010. "Money Laundering Using Trust and Company Service Providers." Paris: Financial Action Task Force - Organization for Economic Cooperation and Development. <http://www.fatf-gafi.org/media/fatf/documents/reports/Money%20Laundering%20Using%20Trust%20and%20Company%20Service%20Providers..pdf>.
- . 2012. "International Standards on Combating Money Laundering and the Financing of Terrorism & Proliferation. The FATF Recommendations." Paris, France: Financial Action Task Force - Organization for Economic Cooperation and Development. http://www.fatf-gafi.org/media/fatf/documents/recommendations/pdfs/FATF_Recommendations.pdf.
- . 2014. "FATF Guidelines on Transparency and Beneficial Ownership." Paris: Financial Action Task Force - Organization for Economic Cooperation and Development. <http://www.fatf-gafi.org/media/fatf/documents/reports/Guidance-transparency-beneficial-ownership.pdf>.
- . 2017. "Topic: High-Risk and Non-Cooperative Jurisdictions." Financial Action Task Force. [http://www.fatf-gafi.org/publications/high-riskandnon-cooperativejurisdictions/?hf=10&b=0&s=desc\(fatf_releasedate\)](http://www.fatf-gafi.org/publications/high-riskandnon-cooperativejurisdictions/?hf=10&b=0&s=desc(fatf_releasedate)).
- . 2019. "Procedures for the FATF Fourth Round of AML/CFT Mutual Evaluations." Financial Action Task Force. <https://www.fatf-gafi.org/publications/mutualevaluations/documents/4th-round-procedures.html>.
- Fazekas, Mihály, István János Tóth, and Lawrence Peter King. 2013. "Anatomy of Grand Corruption: A Composite Corruption Risk Index Based on Objective Data." *Social Science Research Network* 2: 1–46.

- Ferwerda, Joras, and R. Edward Kleemans. 2018. "Estimating Money Laundering Risks: An Application to Business Sectors in the Netherlands." *European Journal of Criminal Policy and Research*.
- Garcia-Bernardo, Javier, Jan Fichtner, Frank W. Takes, and Eelke Heemskerk. 2017. "Uncovering Offshore Financial Centers: Conduits and Sinks in the Global Corporate Ownership Network | Scientific Reports." *Nature Scientific Reports* volume 7, Article number: 6246 (2017). <https://www.nature.com/articles/s41598-017-06322-9>.
- Halliday, Terence, Michael Levi, and Peter Reuter. 2014. "Global Surveillance of Dirty Money: Assessing Assessments of Regimes To Control Money-Laundering and Combat the Financing of Terrorism." Chicago: American Bar Foundation. http://www.lexglobal.org/files/Report_Global%20Surveillance%20of%20Dirty%20Money%201.30.2014.pdf.
- Hangacova, Natalia, and Tomas Stremy. 2018. "Value Added Tax and Carousel Fraud Schemes in the European Union and the Slovak Republik." *European Journal of Crime, Criminal Law and Criminal Justice*. <https://doi.org/10.1163/15718174-02602005>.
- He, H., and E. A. Garcia. 2009. "Learning from Imbalanced Data." *IEEE Transactions on Knowledge and Data Engineering*, September 2009.
- HM Revenue & Customs. 2010. "Anti-Money Laundering Guidance for Trust or Company Service Providers." London: HM Revenue & Customs. https://www.gov.uk/government/uploads/system/uploads/attachment_data/file/372271/mlr8_tcsp.pdf.
- IADB, and OECD. 2019. *A Beneficial Ownership Implementation Toolkit*. Inter-American Development Bank and OECD. <https://doi.org/10.18235/00017111>.
- James, Gareth, Daniela Witten, Trevor Hastie, and Robert Tibshirani. 2014. *An Introduction to Statistical Learning: With Applications in R*. Springer Publishing Company, Incorporated.
- Knobel, Andres. 2019a. "Beneficial Ownership in the Investment Industry: A Strategy to Roll Back Anonymous Capital." *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.3470358>.
- . 2019b. "Beneficial Ownership Verification: Ensuring the Truthfulness and Accuracy of Registered Ownership Information." TJN with the support from Financial Transparency Coalition.
- Knobel, Andres, Moran Harari, and Markus Meinzer. 2018. "The State of Play of Beneficial Ownership Registration: A Visual Overview." SSRN Scholarly Paper ID 3204532. Rochester, NY: Social Science Research Network. <https://papers.ssrn.com/abstract=3204532>.
- Levi, Michael, Peter Reuter, and Terence Halliday. 2018. "Can the AML System Be Evaluated without Better Data?" *Crime, Law and Social Change* 69 (2): 307–28. <https://doi.org/10.1007/s10611-017-9757-4>.
- Netherlands Chamber of Commerce - KVK. 2020. "Foundation - Stichting." Business.Gov.Nl. 2020. <https://business.gov.nl/starting-your-business/choosing-a-business-structure/foundation/>.
- OECD. 2001. "Behind the Corporate Veil: Using Corporate Entities for Illicit Purposes." Paris: Organisation for Economic Co-operation and Development. <http://www.oecd.org/daf/ca/behindthecorporateveilusingcorporateentitiesforillicitpurposes.htm>.
- . 2019. "Peer Review Report on the Exchange of Intormation on Request." Global Forum on Transparency and Exchange of Information for Tax Purposes - Organisation for Economic Co-operation and Development. <https://www.oecd.org/tax/transparency/global-forum-on-transparency-and-exchange-of-information-for-tax-purposes-the-netherlands-2019-second-round-fdce8e7f-en.htm>.

- Pedregosa, Fabian, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, et al. 2011. "Scikit-Learn: Machine Learning in Python." *Journal of Machine Learning Research*, 2011.
- Riccardi, Michele. 2020. *Beyond Blacklists: An Alternative Approach to Identifying Countries at High-Risk of Money Laundering and Illicit Financial Flows - Working Paper*.
- Riccardi, Michele, and Riccardo Milani. 2018. "Opacity of Business Ownership and the Risk of Money Laundering." In *The Janus Faces of Cross-Border Crime in Europe*, Eleven International publishing.
- Riccardi, Michele, Riccardo Milani, and Diana Camerini. 2018. "Assessing Money Laundering Risk across Regions. An Application in Italy." *European Journal of Criminal Policy and Research*.
- Riccardi, Michele, and Ernesto U. Savona, eds. 2013. *Final Report of Project BOWNET - Identifying the Beneficial Owner of Legal Entities in the Fight against Money Laundering Networks*. Trento: Transcrime - Università degli Studi di Trento. <http://www.bownet.eu/materials/BOWNET.pdf>.
- Steinko, Armando Fernández. 2012. "Financial Channels of Money Laundering in Spain." *British Journal of Criminology* 52 (5): 908–31.
- Tavares, Rui. 2013. "Relationship between Money Laundering, Tax Evasion and Tax Havens." Thematic Paper on Money Laundering. Bruxelles: European Parliament - Special Committee on Organised Crime, Corruption and Money Laundering. http://www.europarl.europa.eu/meetdocs/2009_2014/documents/crim/dv/tavares_ml/_tavares_ml_en.pdf.
- Tax Justice Network. 2015. "Financial Secrecy Index 2015 - Final Results." Tax Justice Network. <http://www.financialsecrecyindex.com/PDF/FSI-Rankings-2015.pdf>.
- . 2018. "Financial Secrecy Index 2018 - Methodology." <http://www.financialsecrecyindex.com/PDF/FSI-Methodology.pdf>.
- . n.d. "Regulating Complex Ownership Chains – a Call for Experts for Closed Brainstorming Round-Table Session." Accessed November 11, 2020. <https://www.taxjustice.net/2020/09/17/regulating-complex-ownership-chains-a-call-for-experts-for-closed-brainstorming-round-table-session/>, <https://www.taxjustice.net/2020/09/17/regulating-complex-ownership-chains-a-call-for-experts-for-closed-brainstorming-round-table-session/>.
- Transcrime, ed. 2018. *Mapping the Risk of Serious and Organised Crime Infiltration in Europe - Final Report of the MORE Project*. Milano: Università Cattolica del Sacro Cuore. www.transcrime.it/more.