

Ash, Elliott; Galletta, Sergio; Giommoni, Tommaso

**Working Paper**

## A Machine Learning Approach to Analyze and Support Anti-Corruption Policy

CESifo Working Paper, No. 9015

**Provided in Cooperation with:**

Ifo Institute – Leibniz Institute for Economic Research at the University of Munich

*Suggested Citation:* Ash, Elliott; Galletta, Sergio; Giommoni, Tommaso (2021) : A Machine Learning Approach to Analyze and Support Anti-Corruption Policy, CESifo Working Paper, No. 9015, Center for Economic Studies and Ifo Institute (CESifo), Munich

This Version is available at:

<https://hdl.handle.net/10419/235385>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*

## A Machine Learning Approach to Analyze and Support Anti- Corruption Policy

*Elliott Ash, Sergio Galletta, Tommaso Giommoni*

## **Impressum:**

CESifo Working Papers

ISSN 2364-1428 (electronic version)

Publisher and distributor: Munich Society for the Promotion of Economic Research - CESifo GmbH

The international platform of Ludwigs-Maximilians University's Center for Economic Studies and the ifo Institute

Poschingerstr. 5, 81679 Munich, Germany

Telephone +49 (0)89 2180-2740, Telefax +49 (0)89 2180-17845, email [office@cesifo.de](mailto:office@cesifo.de)

Editor: Clemens Fuest

<https://www.cesifo.org/en/wp>

An electronic version of the paper may be downloaded

- from the SSRN website: [www.SSRN.com](http://www.SSRN.com)
- from the RePEc website: [www.RePEc.org](http://www.RePEc.org)
- from the CESifo website: <https://www.cesifo.org/en/wp>

# A Machine Learning Approach to Analyze and Support Anti-Corruption Policy

## Abstract

Can machine learning support better governance? In the context of Brazilian municipalities, 2001-2012, we have access to detailed accounts of local budgets and audit data on the associated fiscal corruption. Using the budget variables as predictors, we train a tree-based gradient-boosted classifier to predict the presence of corruption in held-out test data. The trained model, when applied to new data, provides a prediction-based measure of corruption that can be used for new empirical analysis or to support policy responses. We validate the empirical usefulness of this measure by replicating and extending some previous empirical evidence on corruption issues in Brazil. We then explore how the predictions can be used to support policies toward corruption. Our policy simulations show that, relative to the status quo policy of random audits, a targeted policy guided by the machine predictions could detect almost twice as many corrupt municipalities for the same audit rate. Similar gains can be achieved for a politically neutral targeting policy that equalizes audit rates across political parties.

JEL-Codes: D730, E620, K140, K420.

Keywords: algorithmic decision-making, corruption policy, local public finance.

*Elliott Ash*  
*ETH Zürich / Switzerland*  
*ashe@ethz.ch*

*Sergio Galletta*  
*University of Bergamo / Italy*  
*sergio.galletta@unibg.it*

*Tommaso Giommoni*  
*ETH Zürich / Switzerland*  
*giommoni@ethz.ch*

We benefited from comments by participants at seminars at the University of St. Gallen, ETH Zürich, University of Zürich, the “Advances in Economics Winter Symposium” (Bergamo), Bonn Law and Economics Workshop, and the Online Workshop in Computational Analysis of Law. We are grateful to Francisco Cavalcanti for sharing essential data for conducting the analysis. David Cai, Matteo Pinna, and Angelica Serrano provided excellent research assistance. Thanks to Daniel Bjorkegren, Aniket Kesari, Himabindu Lakkaraju, Michael Livermore, and Julian Nyarko for comments on previous drafts.

*This version: April 2021 (First version: April 2020).*

## 1. Introduction

A large body of anecdotal and empirical evidence speaks to the deep and negative impacts of corruption. According to recent United Nations statistics, for example, international corruption costs the global economy over 3.6 trillion USD annually. On a more micro level, social scientists have demonstrated that foul play by government actors does real harm to the average citizen. These harms lead to responses in politics and political participation (Ferraz and Finan, 2008; Chong et al., 2015), undermine trust toward institutions (Morris and Klesner, 2010), and have additional side effects on the economy (Lagaras et al., 2017).

Accordingly, researchers continue to seek a scientific understanding of corruption. Broadly speaking, the previous research has identified two important factors. First, electoral incentives play a crucial role in discouraging misbehavior by officials (Ferraz and Finan, 2008; Winters and Weitz-Shapiro, 2013). Second, an effective judicial system to prosecute offenders and enforce the law may be necessary to deter corrupt actions (Becker, 1968; Djankov et al., 2003). Despite this impressive progress in understanding the causes and consequences of corruption, a major impediment to further research is the relative lack of data on corruption. Corrupt actors have strong incentives to conceal their actions, and therefore measurements of corruption traditionally come from costly government auditing programs.

The difficulties facing corruption research also apply to anti-corruption policy efforts. Even with accountable politicians and with well-functioning courts, anti-corruption policies are still often frustrated by the costs of detecting corruption in the first place. Hence, although several countries have introduced monitoring programs to detect wrongdoing, these are typically limited to a relatively small subset of public offices. Overall, producing more data on corruption has high social value in terms of social science research, policy experimentation, and strengthening enforcement.

This paper aims to address the problem of undetected corruption using tools from machine learning. The core of our idea is to exploit the fact that corruption, by its nature, is related to how politicians and public officials manage public resources (Mauro, 1998). Our analysis focuses on corruption in local public finances in Brazilian municipalities. We start with a ground-truth measure of detected corruption, identified and quantified by professional government auditors (Ferraz and Finan, 2008; Brollo et al., 2013). We link this corruption outcome with a rich historical account of local public budgets (with

information on 797 fiscal categories).

We use machine learning to predict corruption from the features of the budget accounts. We implement a gradient boosted classifier consisting of an ensemble of decision trees, typically used to identify patterns in high-dimensional datasets. Using only municipal budget characteristics, the classifier can detect the existence and predict the intensity of corruption with high accuracy in held-out (unseen) data. In the best model, we get an accuracy of 72% and an AUC of 0.77, far better than guessing the modal category or prediction using linear models.<sup>1</sup> We show that the model accurately ranks municipalities by probability of corruption, reproducing the distribution of corruption in held-out data. In a dataset of municipalities that were audited twice, the model can predict within-municipality changes in corruption over time.

To better understand the black box ensemble, we use model explanation techniques to identify the pivotal budget factors that the model attends to when making its predictions. The resulting feature importance scores identify intuitive factors in the budget that are anecdotally related to corruption. In a more quantitative validation, we show using text analysis of the written audit reports that the pivotal budget factors also tend to be mentioned in the reports.

We then take the trained and validated model to the rest of the unlabeled budget data, and form a synthetic measure of corruption for all municipalities and years. To demonstrate the empirical applicability of the method, we use the predicted corruption measure to replicate previous causal results on local corruption in Brazil. First, we replicate the result from [Brollo et al. \(2013\)](#) that a revenue windfall, based on population thresholds, increases corruption. In particular, we can show this result in an untouched sample of municipalities that were never audited by the Brazilian authorities. Normalized coefficient magnitudes are comparable with the estimates obtained by [Brollo et al. \(2013\)](#) using the auditor-produced corruption label as the outcome.

As a second empirical application, we extend the analysis from [Avis et al. \(2018\)](#) and analyze the causal effect of auditing on corruption. Because we have a measure of corruption by year, we can implement an event study analysis. We show that audits reduce corruption in fiscal accounts over the subsequent years, with an average drop of around 2.7% in the probability of malfeasance. Moreover, the effect is especially large

---

<sup>1</sup>As a reference, our classifier’s performance is similar to the model for predicting recidivism from [Kleinberg et al. \(2018\)](#).

for audits that did find corruption, with an average decline of around 18%, about half of the pre-audit mean of 39.5%. In comparison, there is no effect on our measure for audits that did not find corruption.

Besides empirical analysis, our machine predictions for corruption can also be used as an input to policy-making. In the last part of the paper, we investigate the potential of an audit policy guided by predictions on corruption risk. We show that, compared to the status quo policy of random audits, a targeted approach based on predicted corruption would be significantly more efficient in the policy goals of detecting and reducing corruption. According to our policy simulations, a targeted approach would detect 84 percent more corrupt municipalities relative to the random lottery (for the same number of implemented audits). Similarly, by targeting the municipalities at highest risk for corruption, the audit agency could obtain the same number of corruption detections as the random-lottery system but with 45 percent fewer audits, with a corresponding reduction in administrative costs. From a deterrence perspective, it is notable that the annual audit probability, conditional on being corrupt, increases from 3.6 percent under random audits to 6.7 percent under targeted audits.

Finally, we consider the implementation issue that algorithmic targeting could differentially affect the audit rates across political parties. We show using the party affiliations of municipal mayors that this bias turns out to be relevant in our setting, as there is substantial variation across the five main parties in targeting incidence relative to random audits. To address this potential barrier to implementation, we draw on recent developments in algorithmic fairness ([Barocas et al., 2019](#); [Rambachan et al., 2020](#)) and adjust our audit targeting policy to equalize audit rates across parties. We show that such a fair targeting policy can achieve similar gains in policy effectiveness (detecting more corruption) while remaining politically neutral.

Our findings are related to several literatures in economics. First, our paper contributes to the literature studying the relation between corruption and public finance. Many studies emphasize the connection between governmental transfers and public corruption: [Brollo et al. \(2013\)](#) focus on the Brazilian setting, while [De Angelis et al. \(2020\)](#) study the impact of European funds on rent-seeking activity. Another set of papers analyze the extent to which corruption originates from public spending ([Hessami, 2014](#); [Cheol and Mikesell, 2018](#)), and there is evidence that policies that constrain public expenditure may reduce corruption ([Daniele and Giommoni, 2020](#)). Further, other works attend to the link between public procurement and rent-seeking ([Conley and Decarolis,](#)

2016; Coviello and Gagliarducci, 2017). Our results confirm the deep link between public financing and corruption with a focus on the entire budget, instead of single elements, to explain malfeasance. Our approach has the advantage of being general, making it possible to capture the complementary aspects within the budget.

In particular, we add to the existing evidence on the efficacy of auditing programs on corruption in developing countries. Olken (2007) set up an RCT with villages in Indonesia and find that the introduction of the auditing scheme decreased corruption. Bobonis et al. (2016), studying municipalities from Puerto Rico, show that audits effectively reduce corruption and rent-seeking activities by enhancing electoral accountability in the short run, but these effects do not last. Zamboni and Litschig (2018) show in the Brazilian context that increasing the probability of being audited was already effective in reducing corruption. Avis et al. (2018) also study the Brazilian case and find that the implementation of an audit in a specific city reduces future corruption levels in that city. Our event study analysis confirms the latter results, and we are the first to show the dynamics of this effect. Moreover, we find that the effect is particularly strong in cities where corruption is actually detected.<sup>2</sup>

Methodologically, our study adds to the emerging literature in economics applying machine learning techniques to overcome limitations of standard datasets (Athey, 2018). The most established technique in empirical work is to use unsupervised learning to analyze high-dimensional data. For example, Hansen et al. (2018) use Latent Dirichlet Allocation (an unsupervised machine learning algorithm) to measure topics and diversity of discussion in Central Bank committee meeting transcripts. Bandiera et al. (2020) use a similar method to detect CEO behavioral types from their work activity records. Like these papers, we use machine learning to extract relevant dimensions from high-dimensional data. However, we use supervised learning (rather than unsupervised learning) to construct these measurements. This approach is related to several papers in political economy that have used supervised learning to extract measures of partisanship from text, to show (for example) changes in polarization over time or to analyze media influence (Gentzkow and Shapiro, 2010; Ash et al., 2017; Gentzkow et al., 2019; Widmer et al., 2020).

---

<sup>2</sup>Our study also contributes to the body of work on corruption and politics in Brazil. For instance, Ferraz and Finan (2008) show that the disclosure of scandals reduces vote shares for the incumbent. Cavalcanti et al. (2018) emphasize that exposing corrupted incumbents affects the quality of candidates selected by their party to run in the following election.



At the intersection of machine learning and development economics, several papers have applied machine learning methods to detect corruption. The closest paper is [Colonnelli et al. \(2019\)](#), who also predict the results of corruption audits in Brazilian municipalities but focusing on non-budget variables (private sector activity, financial development, and human capital measures). Besides our focus on fiscal factors, the main difference in our paper is to use the measure of corruption for an empirical analysis and policy simulation analysis.<sup>3</sup>

Our use of machine learning to guide auditing is most relevant to the literature on AI-powered policy design ([Kleinberg et al., 2015](#); [Athey, 2018](#); [Knaus et al., 2018](#); [Athey and Wager, 2021](#)). In particular, our approach and results complement those produced by [Kleinberg et al. \(2018\)](#) for criminal recidivism. That paper shows how an algorithm can support the decisions of judges on pre-trial bail release, finding that the algorithm can effectively reduce recidivism by identifying which offenders should be denied bail. Correspondingly, we show that machine learning can support government efforts to identify municipalities with suspicious public budgets, where further investigation is warranted. Other work in this vein has used machine learning to detect higher-quality teachers ([Rockoff et al., 2011](#)), support physician decision-making ([Kleinberg et al., 2015](#); [Mullainathan and Obermeyer, 2019](#)), identify restaurants for targeted health inspections ([Kang et al., 2013](#); [Glaeser et al., 2016](#)), allocate tax rebates and tax audits ([Andini et al., 2018](#); [Battiston et al., 2020](#)), assign refugees to their economically optimal locations ([Bansak et al., 2018](#)), demarcate areas of the Amazon for protection against deforestation ([Assunção et al., 2019](#)), or identify individuals who are most responsive to marketing nudges ([Hitsch and Misra, 2018](#); [Knittel and Stolper, 2019](#)). Besides the new setting (corruption policy), we expand on this work in several methodological directions. First, we use model explanation to validate how the model makes its predictions. Second, we validate the empirical relevance of our machine predictions by showing that they respond appropriately as outcomes in causal regressions. Third, we adopt methods from algorithmic fairness (e.g. [Rambachan et al., 2020](#); [Kasy and Abebe, 2020](#)) to address potential political biases in the targeted audits.

The paper is organized as follows. In [Section 2](#) we present the institutional setting

---

<sup>3</sup>In addition, [López-Iturriaga and Sanz \(2018\)](#) predict the presence of a corruption case each year in 52 Spanish provinces. More at the micro level, [Gallego et al. \(2018\)](#) predict corruption investigations associated with a sample of 2 million public contracts in Colombia.

and the data. Section 3 describes the prediction procedures and model performance results. In Section 4 we provide the estimation strategy and the results of the empirical applications, while Section 5 reports our policy simulations for guided audits. Section 6 concludes.

## 2. Institutional Background and Data Sources

### 2.1. Local Government and Budgets

Brazil has a decentralized governance structure composed of 26 states and 5563 municipalities. At the municipal level, the central political authorities are the mayor (*prefeito*) and the city council (*Câmara de Vereadores*), which are directly elected by citizens every 4 years. Starting from the 1980s, local governments have enjoyed substantial autonomy in public budgeting decisions. They have primary responsibility for the provision of health and education services and municipal transportation and infrastructure. For the most part, these services are funded by upper-level jurisdictions via intergovernmental transfers. Yet, the mayor has autonomy in setting the tax rate for important local taxes, e.g., taxes on buildings and lands (*Imposto sobre a Propriedade Predial e Territorial Urbana* - IPTU), as well as sales taxes on services (*Imposto sobre Serviços*).

We collected the annual budget of all Brazilian municipalities for 2001 through 2012. Building on the previous local public finance literature, we gather detailed information about the categories of expenditure, revenue, active positions (assets), and passive positions (liabilities). These data are publicly available in the Finance Ministry’s online database.<sup>4</sup> We downloaded the datasets for each year and cleaned the variables to make them comparable across years.

In the period of our analysis, the budgets were composed of a large number of different categories for each of the four macro-categories. In total, we have 797 accounting variables from the original data source. The expenditure section has the most components, while the passive section has the fewest. Over time, there is an increasing level of detail about the use and sources of local governments’ revenue as the budget adapts to changes in legislation. There is some missingness, as not all categories are reported for each year and municipality. Appendix Table A1 reports the number of categories for

---

<sup>4</sup><https://www.tesourotransparente.gov.br>

each section of the balance sheet for each year of data. We impute missing values using the mean value of the associated variable.

## 2.2. *Anti-corruption policy in Brazil*

In 2003, the Brazilian government introduced new policies to reduce corruption. In particular, the policymakers behind this agenda were concerned about misuse of federally transferred funds by local authorities. Thus, a cornerstone of the reform was a system of random audits, in which municipalities are randomly selected to have their fiscal accounts audited for corruption.

The government invested significant planning and resources in these inspections. In particular, random assignment of audits was implemented to ensure fairness in their allocation. In a given audit round, of which there are around four per year, between 50 and 60 municipalities are chosen. Separate lotteries are run for each state (meaning some states getting slightly more lotteries per municipality than others), and cities with more than 500,000 inhabitants are excluded. Otherwise, the audits are exogenously assigned.

The audits are implemented by officials from *Controladoria Geral da União* (CGU), an independent federal public agency. Every selected municipality is visited by 10 to 15 auditors. Their inspections focus on a list of randomly selected items provided by the CGU from the sample of federal transfers the municipality received in the previous 3-4 years. They usually spend a couple of weeks in municipal offices collecting information to identify potential mismanagement in the use of public funds. The auditors summarize the presence of irregularities in reports made available to the public within a few months of the inspection. These audit reports provide detailed information that can be used to create measures of municipal-level corruption (Ferraz and Finan, 2008; Brollo et al., 2013; Zamboni and Litschig, 2018). We use the corruption measures provided by Brollo et al. (2013). These data include several measures for all 1,481 municipalities audited in the first 29 lotteries of the anti-corruption program (i.e., audits from 2003 to 2009). Focusing on a particular mayor’s term of office, they compute the share of corrupted resources (i.e., the ratio between the total amount of funds involved in the detected violation and the total amount audited). Our analysis focuses on a binary variable identifying the presence of what the authors call *narrow corruption*, which is restricted to severe

Table 1: Summary Statistics

Variable	Mean	Std. Dev.	Min	Max	N
<b><i>True corruption (term)</i></b>					
Main Labels from <a href="#">Brollo et al. (2013)</a>	0.424	0.494	0.000	1.000	2087
Alternative Labels from <a href="#">Avis et al. (2018)</a>	0.238	0.426	0.000	1.000	1604
<b><i>Budget categories (year)</i></b>					
Tax on agricultural territorial property (ITR)	4.2	21.2	-0.0	1414.2	64933
Spending in agriculture	31.2	85.7	0.0	10108.5	64933
Spending in transportation	69.6	127.0	0.0	10155.0	64933
Tax on export of industrialized products (IPI)	4.7	8.3	-1.1	431.0	64933
Budget Surplus/Deficit	41.0	3339.2	-3743.5	650900.8	64933
Cash	3.5	35.0	-1607.5	5017.8	64933
Tax on real estate transactions (ITB)	9.7	18.6	-0.0	917.0	64933
Taxes	8.9	15.9	0.0	781.5	64933
Deposits	21.9	72.7	-468.1	12335.8	64933
Motor vehicle property tax (IPVA)	19.1	27.3	0.0	2120.0	64933
<b><i>Municipal characteristics</i></b>					
Mean income	593.0	319.8	29.8	3062.5	64933
Agriculture (% employed)	16.9	10.1	0.0	72.3	64933
Industry (% employed)	4.2	4.2	0.0	37.5	64933
Commerce (% employed)	7.5	3.6	0.3	27.8	64933
Transport (% employed)	1.2	0.7	0.0	5.9	64933
Service (% employed)	6.8	2.7	0.3	19.3	64933
Public administration (% employed)	2.1	1.2	0.1	16.1	64933
Employed population	38.4	8.5	9.7	79.8	64933
Graduated people	1.2	1.3	0.0	16.5	64933
Poor population	10.0	8.1	0.3	54.4	64933
Gini coefficient	0.6	0.1	0.3	0.9	64933

*Notes:* *Main Labels* from [Brollo et al. \(2013\)](#) captures the binary variable measuring the presence of corruption according to [Brollo et al. \(2013\)](#) (*narrow corruption* variable). *Alternative Labels* from [Avis et al. \(2018\)](#) captures the binary variable measuring the presence of corruption according to [Avis et al. \(2018\)](#). All budget variables are expressed in per-capita terms. The municipal characteristics are drawn from the 2000 Brazilian census. *Mean income* captures the average income of the working population, the variables *Agriculture*, *Industry*, *Commerce*, *Transport*, *Service* and *Public administration* capture the population employed in a specific sector. *Employed population* measures the fraction of employed population, *Graduated people* is expressed in percentage points and *Poor population* is the fraction of poor population.

irregularities such as illegal procurement, fraud, favoritism, and over-invoicing.<sup>5</sup> On this definition, 42% of audited municipalities at their first audit are found to be corrupt.

For robustness, we have access to an alternative set of corruption labels from [Avis et al. \(2018\)](#). This measure is constructed using a slightly different approach to coding the audit report documents. It is available for a different (but mostly overlapping) set of audits. We find that the two measures are highly correlated (Appendix Figure A3). In Appendix B, we will provide supplementary analysis using this alternative measure of corruption.

<sup>5</sup>In addition, they define a measure of *broad corruption*, which also includes inconsistencies that could be linked to government mismanagement, but not intentional misuse. This concept of corruption is less useful because it is so widespread: 76% of audited municipalities have broad corruption.

### 2.3. Linked Dataset

We join the corruption labels, which are defined at the municipality-term level, with the local budget factors, which are defined at the municipality-year level. The resulting dataset is at the municipality-year level. We then add data on local demographics, on intergovernment transfers, and political party control. Specifically, we add demographics from the 2000 Brazilian Census, including *mean income*, *share of population employed*, *sector of occupation* (agriculture, industry, commerce, transportation, services and public administration), *share with college education*, *poverty rate*, and *Gini Coefficient of income*. Federal-to-municipal revenue transfers data come from the Brazilian National Treasury (*Tesouro Nacional*). Population data from the Brazilian Institute of Geography and Statistics (IBGE). Finally, we collected information about the mayor party affiliations in the 2000, 2004, and 2008 elections. Summary statistics on these variables are reported in Table 1.

## 3. Predicting Corruption from Budget Data

Our goal is to take the information in the municipal budget and learn a prediction function to provide a probability that a given municipality is fiscally corrupt. To that end, this section outlines how we build our dataset and machine learning model to form those predictions. We evaluate and interpret the predictive model, and then apply it to all municipalities in Brazil for use in the subsequent analysis.

### 3.1. Corruption Prediction Dataset

Our data consists of budget predictors and corruption labels. For the budget features  $X$ , we don't undertake any additional pre-processing steps besides imputing missing values with the mean value for the associated variable.<sup>6</sup> The resulting matrix  $X$  of budget factors has 897 columns, corresponding to the budget fields, and rows corresponding to each municipality and year.

The corruption label  $Y \in \{0, 1\}$ , defined at the municipality-year level during terms subject to audit, equals one for years where an audit found narrow corruption, and equals zero for years when the audit did not find narrow corruption. For the machine learning part, any municipality-terms that were not audited have to be excluded because we do

---

<sup>6</sup>We got similar results when experimenting with additional pre-processing steps, including adding missing indicator variables, (standardizing variables, or transforming variables as per capita.

not have any labels. When municipalities were randomly audited more than once, we exclude the second audit from the machine learning dataset.

### 3.2. Machine Learning Approach

We face a binary classification task. We want to learn a conditional expectation function  $Y(X)$  that provides a predicted probability that a municipality is corrupt based on the publicly observed budget features. Economists are already familiar with logit or probit (for example) in this setting. But these classical statistical models do not extrapolate well to new datasets because they tend to over-fit the training sample (e.g. [Hastie et al., 2009](#)). The contribution of machine learning tools, now becoming widespread in economics (e.g. [Belloni et al., 2014](#); [Mullainathan and Spiess, 2017](#); [Athey, 2018](#)), is to address the over-fitting problem and provide robust out-of-sample prediction with high-dimensional datasets.

Researchers and policymakers now have access to a variety of machine learning tools for solving binary classification tasks. For example, one of the baseline models that we will use below is penalized logistic regression. This model is very similar to the binary logit, which learns a set of linear coefficients on  $X$ , sums them, and puts them through a sigmoid transformation to obtain a probability for  $Y$  between zero and one. What is new is a penalty term, which adds an additional cost to the training objective that penalizes larger coefficients. The penalty addresses overfitting and helps the model predict better in held-out test set data. During the training process the strength of the penalty is calibrated in a process called cross-validation, where the training data is split up and the out-of-sample performance of different penalties is evaluated. Then the best model is taken to the unseen test set for a clean evaluation and for any downstream tasks.

A state-of-the-art model for binary classification using high-dimensional tabular datasets is gradient boosted trees ([Friedman, 2001](#); [Hastie et al., 2009](#)).<sup>7</sup> Gradient boosting models consist of an ensemble of decision trees that “vote” on the predicted outcome. Each decision tree iteratively selects informative variables (e.g., property taxes), splits on a value of that variable (e.g.,  $x > 100$ ), branches off for additional splitting, and so on, until reaching a terminal node and an associated prediction ( $\hat{Y} = 0$  or  $\hat{Y} = 1$ ). With gradient boosting, additional layers of trees are gradually added during the training process to fit residuals and fix errors in the initial layers. This iterative growth approach

---

<sup>7</sup>This is the same algorithm used by [Kleinberg et al. \(2018\)](#) in predicting criminal recidivism.

tends to perform better than other ensemble methods, such as random forests, which grow trees in parallel.

More specifically, we train a gradient boosted classifier using the implementation from the python package XGBoost (Chen and Guestrin, 2016). Feurer et al. (2018) systematically compared XGBoost to many other classifiers, including a sophisticated automated ML system, and found that XGBoost consistently performed best on our type of machine learning task. We used cross validation grid search to tune hyperparameters, which include the learning rate, L1 and L2 regularization penalties on the learned parameter weights, the max depth of the constituent decision trees, and an additional regularization constraint specifying a minimum threshold for the size of decision tree terminal nodes. Appendix Table A2 shows the selected values for these hyperparameters across each of five different training folds.

In the next subsection on model performance, we compare XGBoost to a number of baselines. First, as the weakest baseline, we guess the modal category (not corrupt). Second, we train ordinary least squares (OLS), or non-penalized linear regression, dropping multi-collinear predictors. Third, LASSO, perhaps the most familiar machine learning model to economists (e.g. Belloni et al., 2014), is a linear regression model but adds an L1 penalty that penalizes larger coefficients and outputs a sparse model. For both OLS and LASSO, the predicted probabilities for  $Y$  might be below zero or above one, but a decision threshold of 0.5 is used for assigning a predicted label. Finally, as the strongest alternative baseline, we use penalized logistic regression, a linear classifier with a sigmoid transformation and elastic net penalty (that is, both a LASSO (L1) and a ridge (L2) penalty). For LASSO and Logistic, the penalty is selected by cross-validation grid search in the training set. All three of these linear baselines are implemented using the stochastic gradient descent learners from the python package scikit-learn (Pedregosa et al., 2011).

We train and evaluate models using nested cross-validation, which works as follows: First, we randomly split the sample of audited municipalities into five different sets. Next, we train five separate models using each time four different subsets (80% of the sample) and take the tuned models to get performance metrics in the test set (the remaining 20 % of the sample). Each time, we tune the hyperparameters in the training set using five-fold cross-validation. In each fold, early stopping is used (with patience of ten training epochs) to stop training when the model begins to over-fit the training set. Appendix Table A2 shows that the resulting forests consist of between 46 and 72 trees,

Table 2: Machine Learning Metrics for Predicting Corruption

	OLS (1)	Lasso (2)	Logistic (3)	XGBoost (4)
Accuracy	0.476 (0.022)	0.474 (0.022)	0.560 (0.022)	0.723 (0.012)
AUC-ROC	0.487 (0.016)	0.507 (0.012)	0.568 (0.016)	0.777 (0.013)
F1	0.685 (0.031)	0.538 (0.050)	0.545 (0.054)	0.632 (0.018)

*Notes:* Columns report the mean and standard error (in parentheses) for the indicated performance metrics (by row) across the five model runs, produced using separate training-set folds. Columns indicate the machine learning model used.

each with up to 10 variable splits before a terminal node.

The nested approach provides five sets of predictions for each model. In the model evaluation section, we have five sets of test-set evaluation metrics. We report the mean and standard error across these five models. In the downstream tasks, and in particular the policy simulation, we will use the multiple predictions to assess the importance of sampling variability in the predictions.

### 3.3. Model Performance

We evaluate our set of models by their scores on a set of standard classification metrics in the held-out test data. These metrics, reported by row in Table 2 Panel A, describe how well a model trained on budget accounts can replicate the auditing agency’s judgments about fiscal corruption. First, the most straightforward metric is accuracy, which gives the proportion of test-set observations for which the machine-predicted label matches the true label. A naive guessing model that chooses the modal category (not corrupt) would obtain  $\text{accuracy} = 0.58$ . Second, we report AUC-ROC (area under the receiver operator characteristic curve), another standard metric in binary classification. AUC-ROC, which takes values between 0.5 (random guessing) and 1.0 (perfect accuracy), can be interpreted as the probability that a randomly sampled corrupt municipality is ranked more highly by predicted probability of corruption than a randomly sampled non-corrupt municipality. Third, we report F1 for the corrupt class, defined as the harmonic mean of precision (proportion true corrupt within the set predicted cor-



rupt) and recall (proportion predicted corrupt within the set true corrupt). F1, ranging from 0.0 (guessing the modal category) and 1.0 (perfect accuracy for the corrupt class), penalizes both false positives and false negatives.

Across the columns of Table 2, we compare the predictive performance of our preferred model, XGBoost, to a number of other baselines. In each table cell, we report the average test-set performance metric across the five cross-validated model runs, with the standard error of the mean in parentheses. In the text, we report the minimum and maximum metric values across the folds.

In the rightmost Column 4, we report the metrics from our preferred model, the gradient boosting classifier. The average test set accuracy for the predictions across five nested folds is 0.723, with the minimum accuracy being 0.692 and the maximum 0.755. For AUC-ROC, the average is 0.777 (with min = 0.743 and max = 0.809), while for F1, the average is 0.632 (min = 0.584, max = 0.675). To help contextualize these numbers, the model is similar in its performance to the one used by [Kleinberg et al. \(2018\)](#) to predict criminal recidivism by defendants in pre-trial bail proceedings.<sup>8</sup> For comparison, we form predictions using OLS (Column 1), LASSO (Column 2), and Logistic Regression (Column 3). The predictions provided by OLS or LASSO are barely better than random guessing. Logistic regression is somewhat better but still much worse than XGBoost. This difference in performance suggests a nonlinear, interactional relationship among the predictors that the tree ensemble is better able to learn.

Appendices A and B report additional evaluations of the prediction task. First, to help visualize the distribution of predictions, Appendix Table A3 shows the confusion matrices for the test-set predictions. For XGBoost, we can see good precision and recall across categories. The confusion matrices for OLS, LASSO, and logistic regression show that the linear models tend to produce many false positives (not-corrupt municipalities are often labeled as corrupt).

Second, we focus on the municipalities that have been audited twice and see if our prediction model can reproduce within-municipality changes in corruption over time. To that end, we regress the change in true corruption against the change in predicted corruption, adjusting for audit year fixed effects and demographic characteristics. Appendix Figure A1 shows that there is a significant positive effect in this regression. This within-municipality validation is important for the usefulness of our measure in empirical

---

<sup>8</sup>[Kleinberg et al. \(2018\)](#) report an AUC-ROC of 0.707 for their best-performing model.

tasks, where one would like to be able to examine changes in corruption over time.

Third, in Appendix Table A10 we report performance metrics with an alternative sampling approach and with an alternative corruption label. In Columns 1-3, we apply random splits between training and test set by municipality, instead of by municipality-year, which allows us to compare the model performance using budget factors, fixed demographic factors, or both. The alternative sampling approach obtains comparable performance to our baseline model and shows that a model trained using just demographic information is less accurate than a model using budget information. In Columns 4-7, we show the model performance for the alternative corruption label from Avis et al. (2018). We see that our XGBoost model is even more accurate in predicting the alternative label (AUC-ROC = 0.903, s.e. = 0.009).

### 3.4. Interpreting the Predictions

Gradient boosted machines, like all ensembles and other sophisticated machine learning algorithms, are black boxes. At the end of model training, we have a dense forest of decision trees. With 797 variable being input into those trees, and hundreds of splits within the forest, it is difficult to tell how the model is making its predictions. In this subsection we use model explanation methods to better understand how the model works.

The previous applied machine learning literature has discovered an advantage of gradient boosted machines that compensates for their basic lack of interpretability (e.g. Hastie et al., 2009). One can rank the input variables by their *feature importance*, computed as the number of times the model “uses” that variable in the sense that one of the constituent decision trees splits on it. Note that these features could be either positively or negatively correlated with the predicted corruption. The ranking is more informative than seeing which predictors are correlated with corruption, because they could be important through a non-linear relation, or through interactions with other variables. Moreover, the important features can be seen as *pivotal* in the sense that they are the most useful variables for predicting the outcome, even among clusters of highly collinear predictors.

Here we use the feature importance ranking to get some insight into how our corruption detection model makes its predictions. After model training, we have feature importance scores for each of the five cross-validated models. We average the scores across folds and then rank the most important features.

From the feature importance scores, we learn immediately that our dataset contains

Table 3: Most important budget features for Corruption Prediction

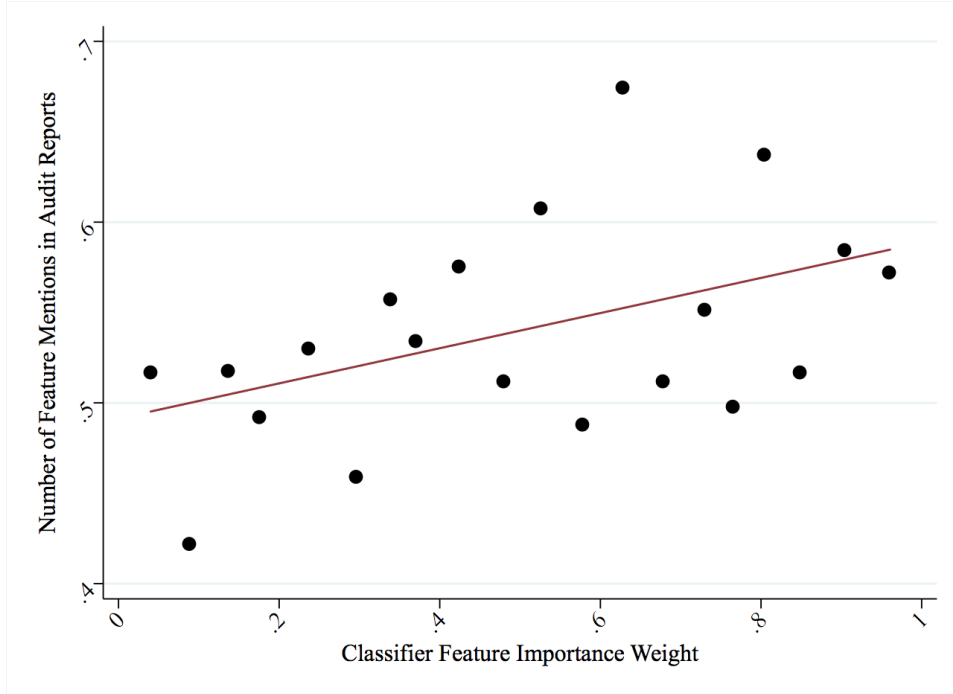
N.	Category	Macro Category	Weight
1	Tax on agricultural territorial property (ITR) (compartecipation)	Revenue	103
2	Spending in agriculture	Expenditure	103
3	Spending in transportation	Expenditure	95
4	Tax on export of industrialized products (IPI) (compartecipation)	Revenue	92
5	Budget Surplus/Deficit		84
6	Cash	Assets	84
7	Tax on real estate transactions (ITB)	Revenue	80
8	Taxes	Revenue	79
9	Deposits	Assets	76
10	Motor vehicle property tax (IPVA) (compartecipation)	Revenue	74
11	Income Tax (IRRF)	Revenue	73
12	Transfers for the health system	Revenue	71
13	Civil servant per diems	Expenditure	70
14	Spending for legislative procedure	Expenditure	68
15	Revenue from assets	Revenue	67

*Notes:* List of the most important features. Metrics rank the features (budget components) by how often they are included in a decision tree contained in the ensemble classifier, averaged across the five training folds.

many noise predictors. Out of the 797 variables input to the ensemble, only 446 are used at all and 351 are ignored by the ensemble across all five folds. Within the set of useful variables, we show the 15 predictors with the highest feature importance scores in Table 3.

The most frequent categories identified as relevant to corruption are those related to expenditures and taxes. On expenditures, we see corruption-related spending for agriculture (2), public services in transportation (3), coherently with [Hessami \(2014\)](#), and legislative actions of local government (14). Other specific signals come from the arbitrary use of public funds in categories that are perhaps more difficult to monitor, for instance, civil servant per diems (13). Many variables are also included on the taxes side ([Liu and Mikesell, 2019](#)), with the model especially attending to income tax (11) and different types of property taxes. These include property taxes on agricultural land (1), on motor vehicles (10), and on transfers of real estate ownership (7). The latter variable is related to the construction sector, traditionally associated to corruption ([Kyriacou et al., 2015](#)). Some of these categories refer to national tax revenues that are transferred to municipalities, as in the case of taxes from export of industrialized products (4). Other non-tax revenues that made it to the list include transfers from other government levels to fund the public health system (12), as studied in [Machoski and de Araujo \(2020\)](#), and revenues generated from municipal assets (e.g., real estate) (15). In terms of assets, we find liquid assets, such as cash (6) and the more general classification of deposits (9).

Figure 1: Model-Predicted Feature Importance and Mentions in Audit Report Texts



*Notes:* Binscatter diagram for frequency that a budget feature appears in the municipal audit reports (vertical axis) against binned feature importance weights for each feature (horizontal axis). Pearson’s correlation is 0.13. The regression coefficient is 0.097 with  $p = .09$  (robust standard errors).

Finally, the model attends to the presence of a budget deficit/surplus, in line with [Liu et al. \(2017\)](#), showing the link between public corruption and debt.

While the model’s identification of these important budget features is consistent with some work from the literature, listing these examples is somewhat ad hoc. To see whether the feature importance scores can validate our model more quantitatively, we would like to know whether these pivotal features were actually identified as related to corruption by the auditors in Brazil. To do that, we look for mentions of these items in the best available place – the text of the published audit reports.

To this end, we downloaded the full library of audit reports for our time period as PDF files from the agency web site. The PDFs were in machine-readable Portuguese and therefore straightforward to extract as plain text. We performed mild cleaning the language, namely removing punctuation and capitalization. The same was done for our list of budget accounting variable names. Finally, we counted the total mentions of each budget feature in the corpus of reports.

We then produced a dataset at the prediction variable level, containing the percentile

rank in the model feature importance score and the percentile rank in the audit-report mention count. Figure 1 plots the audit mention percentiles against the feature importance percentiles. We see a clear positive relationship that is statistically significant in a univariate regression ( $p = .09$ ). Our classifier, trained on the budget accounts with just corruption labels, identifies as important the same budget features that tend to be mentioned in the audit report documents. These validation results support the view that our measure captures activities that are indeed related to corruption.

### 3.5. *Measuring Corruption in Non-Audited Municipalities*

An essential contribution of our approach is to measure corruption for all Brazilian municipalities and all years from 2001 to 2012. Using the trained models, we form five predicted corruption probabilities for all observations based on the budget data. In Figure 2 we provide a visualization of the difference between the sample of only audited municipalities (Panel a) and the sample of municipalities that we can analyze when using our predicted measure of corruption (Panel b). The map illustrates quite clearly the additional information produced by the machine learning method. With the machine predictions, we can then analyze corruption in municipalities (and years) regardless of whether they have been audited.

## 4. Empirical Applications

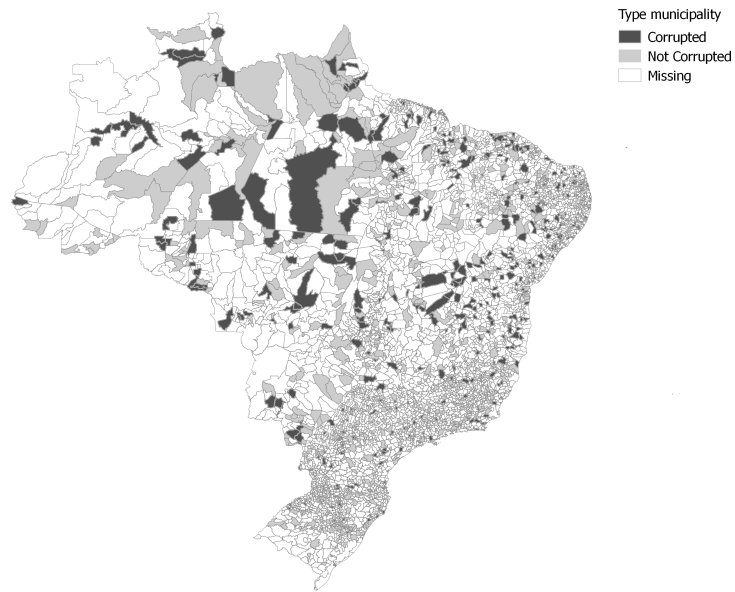
This section replicates and extends existing evidence from the literature on corruption in Brazil. This exercise has two purposes. On the one hand, it provides checks on the internal validity of our synthetic measure of corruption – that is, we can check whether it responds to causal treatments the same way as auditor-measured corruption. On the other hand, we extend previous results by taking advantage of the larger sample of municipalities and the time variation of our corruption measure.

### 4.1. *Revenue Shocks and Corruption*

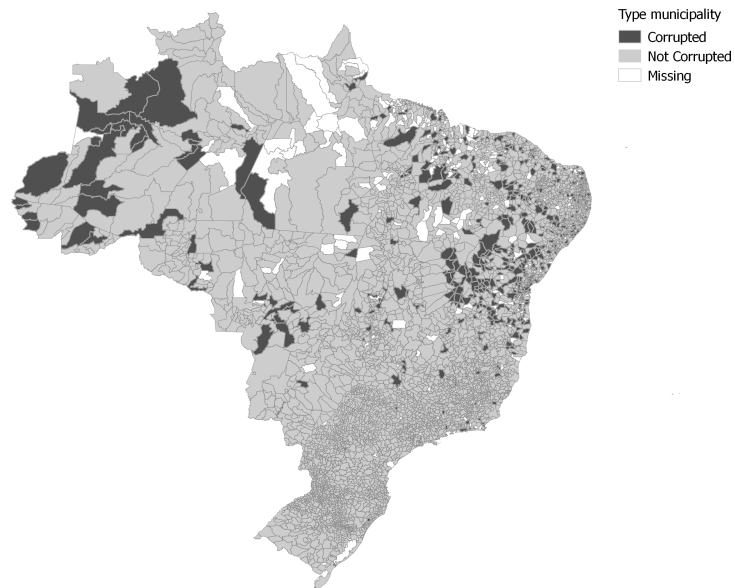
As a first analysis, we use the new synthetic measure of corruption to analyze the effect of revenue shocks on corruption, replicating and extending the findings by Brollo et al. (2013). This paper studies whether a windfall of public revenues can lead to an increase in rent-seeking by the public administration (as measured by a subsequent surge in corruption). They estimate the impact of federal transfers on the occurrence of corruption as detected by the random audits.

Figure 2: The Geography of (Predicted) Corruption

(a) Actual Corruption



(b) Predicted Corruption



*Notes:* Actual (Panel a) and predicted (Panel b) corruption by municipality, using budgets from 2004. A municipality is predicted to be corrupted if mean prediction is  $>0.5$ .

Brazilian municipalities receive transfers from the states and from the federal government. Federal transfers are the largest single source of municipal revenues (around 40% of the total budget). The amount transferred through this *FPM* program (*Fundo de Participação dos Municípios*) depends on exogenous population thresholds, where municipalities in the same state and in a given population bracket receive the same amount of resources.<sup>9</sup>

More precisely, the amount of revenues received by municipality  $i$  in state  $k$  follows the allocation mechanism:

$$FPM_i^k = \frac{FPM_k \lambda_i}{\sum_{i \in k} \lambda_i}$$

where  $FPM_k$  is the total amount allocated in state  $k$  and  $\lambda_i$  is the municipality-specific coefficient, as shown in Table A4. Due to imperfect compliance, however, the statutorily prescribed transfers do not perfectly determine the amounts actually transferred.<sup>10</sup> Thus, Brollo et al. (2013) use a fuzzy regression discontinuity design methodology, instrumenting actual transfers ( $\tau_i$ ) with theoretical transfers ( $\hat{\tau}_i$ ).

Formally, we have the first stage

$$\tau_i = g(P_i) + \alpha_\tau \hat{\tau}_i + \delta_t + \gamma_p + u_i \quad (1)$$

and reduced form

$$y_i = g(P_i) + \alpha_y \hat{\tau}_i + \delta_t + \gamma_p + \eta_i \quad (2)$$

where  $y_i$  is corruption,  $g(\cdot)$  is a high order polynomial in  $P_i$  (the population of city  $i$ ),  $\delta_t$  contains term fixed effects,  $\gamma_p$  contains state fixed effects, and  $u_i$  and  $\eta_i$  are the error terms. The coefficients  $\alpha_\tau$  and  $\alpha_y$  capture the effects of theoretical transfers on actual transfers and (predicted) corruption, respectively. For the two-stage-least squares analysis, we estimate the second stage

$$y_i = g(P_i) + \beta_y \tau_i + \delta_t + \gamma_p + \epsilon_i \quad (3)$$

---

<sup>9</sup>Appendix Table A4 shows these coefficients and the corresponding population brackets: Following Brollo et al. (2013), we focus on the initial seven brackets and restrict the sample to cities with a population below 50,940. Furthermore, we follow the approach of Brollo et al. (2013) and restrict the sample, for the sake of symmetry, to municipalities from 3,396 below the first threshold to 6,792 above the seventh threshold. This sample represents about 90 percent of Brazilian municipalities.

<sup>10</sup>This imperfect compliance is due to many factors (*e.g.* municipalities splitting, manipulation in population figures). See Brollo et al. (2013).

where theoretical transfers  $\hat{\tau}_i$  are used as an instrument for actual transfers  $\tau_i$  and all other terms are defined as above. The coefficient  $\beta_y$  captures the causal effect of actual transfers on (predicted) corruption. For inference, standard errors are clustered by municipality.<sup>11</sup>

Our data cover the two mayoral terms, January 2001–December 2004 and January 2005–December 2008. While [Brollo et al. \(2013\)](#) focus only on municipalities that received an audit, our dataset allows us to analyze a larger and more representative sample of cities. Therefore, our exercise is also providing a test for the external validity of their results.<sup>12</sup> Appendix table [A5](#) shows the descriptive statistics by population bracket. Brazilian municipalities in our sample receive, on average, \$3.3M BRL (about \$610K USD), while theoretical transfers are somewhat higher at \$3.7M BRL (about \$680K USD). The average level of (predicted) corruption is around 0.5 and its level does not change significantly as we move to larger cities.

Table [4](#) reports the results for the regression analysis. Panel A shows the estimates of the first stage, Equation (1), Panel B shows the reduced-form effects, Equation (2), while Panel C shows the the corresponding two-stage-least-squares estimates. For each panel, we provide the results when including: cities that have received an audit (column 1) similarly to [Brollo et al. \(2013\)](#), all cities (column 2), and cities that have never been audited (column 3).

We find a strong first-stage effect, showing that theoretical transfers positively affect actual transfers, and this is true for all samples considered. In addition, we find positive and significant coefficients when estimating the reduced-form as well as the two-stage-least-squares results. Varying the sample of interest does not significantly alter the size of coefficient and the level of precision is stable. Notably, the magnitude of the standardized reduced-form coefficient is about four-fifths the size of that estimated by [Brollo et al. \(2013\)](#), and our 95% confidence interval contains the original coefficient. Thus even the magnitudes of empirical estimates using machine-learning-measured corruption seem to be comparable to using auditor-measured corruption.

To test the robustness of these empirical results, we conducted a series of checks. First, we replicate the main analysis on four random samples of 1,115 municipalities, the

---

<sup>11</sup>See [Brollo et al. \(2013\)](#) for a detailed discussion and testing of the econometric assumptions in this setting.

<sup>12</sup>For the sake of brevity we only replicate the analysis on the overall effect, omitting the threshold-specific analysis.



Table 4: Replication Analysis: Effect of Revenue Shocks on Corruption

	Audited cities (1)	All cities (2)	Non-audited cities (3)
<i>Panel A. First Stage</i>			
Theoretical transfers	0.6805*** (0.0205)	0.6909*** (0.0233)	0.6996*** (0.0230)
<i>Panel B. Reduced Form</i>			
Theoretical transfers	0.0040*** (0.0009)	0.0041*** (0.0003)	0.0040*** (0.0003)
<i>Panel C. 2SLS</i>			
Actual transfers	0.0058*** (0.0013)	0.0059*** (0.0005)	0.0057*** (0.0005)
N. Observations	1115	5808	4693

*Notes:* Effects of FPM transfers on (predicted) corruption measures. Panel A reports the estimates of the first-stage analysis, the dependent variable is *actual transfers*. Panel B reports the estimates of reduced form analysis, the dependent variable is *predicted corruption*. Panel C reports the estimates of the 2sls estimates, the dependent variable is *predicted corruption* and *actual transfers* is instrumented with *theoretical transfers*. Column headings indicate the sample of municipalities included. All regressions controls for a third-order polynomial in normalized population size, term dummies, and macro-region dummies. Robust standard errors clustered at the municipal level are in parentheses: \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

sample size of the original analysis by [Brollo et al. \(2013\)](#) (Appendix Table A6). The coefficients show some variation, but they are always positive and statistically significant. Second, we show that the instrument is not correlated with the error of the prediction model, defined as the difference between the true corruption level and the predicted one (p-value=0.212). This null is helpful because it suggests that the model errors are not responding to the instrument. That is, the correlated factors besides corruption that are contributing to our prediction are not affected directly by revenue transfer shocks. Thus, using our model predictions as the outcome will still satisfy the exclusion restriction. Third, we show in Appendix Table A7 Column 1 that there is no revenue-shock effect on a corruption prediction formed with a model trained on municipal demographic characteristics (similar to Collonelli et al’s). This placebo test is reassuring because the model trained on demographics does not contain budget information, and means that our model is not forming corruption predictions based on spurious correlations with demographics. Fourth, we formed predictions from our baseline model while permuting randomly the FPM transfer variable, which could be mechanically shifted by the revenue shocks instrument. The effect of revenue shocks is the same (Appendix Table A7 Column 2).

Overall, this replication exercise provides helpful validation for the use of our predicted measure of corruption in contexts where audits provide insufficient data. In addition, we provide additional evidence on the external validity of findings by [Brollo et al. \(2013\)](#).

#### *4.2. Effect of Audits on Corruption*

The next empirical application uses our predicted measure of corruption to analyze the effect of auditing on subsequent corruption in an event study framework. This analysis complements [Avis et al. \(2018\)](#), who explore the same research question using the set of Brazilian municipalities that were (by random draw) audited twice in a cross-sectional setup. With our new measure of predicted corruption, we can overcome the data limitations of Avis et al. and extend their results. First, because of the longitudinal nature of our dataset, we can capture dynamic effects. Second, we can condition our estimates on pre-audit levels of corruption. Third, our effects are identified by a relatively

larger sample of municipalities that got audited only once (rather than twice).<sup>13</sup>

Using the annual corruption prediction  $y_{ist}$  in municipality  $i$  of state  $s$  at year  $t$ , we take a standard event study approach and estimate the within-municipality effects of a (randomly assigned) audit. Let  $D_{ist}^k$  be a dummy variable for  $k$  years before and after an audit. We estimate

$$y_{ist} = \sum_{k=-3, k \neq -1}^5 \beta_k D_{ist}^k + \delta_i + \lambda_t + W_{ist}'\phi + \epsilon_{ist} \quad (4)$$

where we have municipality fixed effects  $\delta_i$ , year fixed effects  $\lambda_t$  and other controls  $W_{ist}$ , which in particular includes dummy variables indicating periods distant from when the audit took place. Because  $k \neq -1$  (the year before the audit), the  $\beta_k$  estimate the dynamic effects relative to the year before the audit. The identifying assumption hinges on randomness in the timing of selection into the audit program. We cluster standard errors by state. The sample includes 1,479 municipalities that have received an audit in the time period under analysis.

We graphically report estimates for Equation (4) in Figure 3 Panel (a), with the numerical estimates reported in Appendix Table A8. We can see that already in the year of the audit ( $k = 0$ ), there is a sharp and statistically significant drop in predicted corruption. This persists over the subsequent years but becomes weaker. Meanwhile, as expected given the random assignment due to the lottery, there is no statistically significant effect in the pre-announcement years.

Panel (b) reports event-study effects for the subsets of audits that find clear corruption (black points) and those that do not find corruption at all (grey points).<sup>14</sup> These trends look quite different. When corruption is discovered (black points), there is a much larger negative effect ranging between -1.7% and -25.8%, which is sizeable if compared with the magnitude of the treatment mean of 55.8%. The effect is persistent across subsequent years. In contrast, when the audit does not find any corruption or irregularities (grey points), there is no effect on corruption. Such effects could consist of an actual

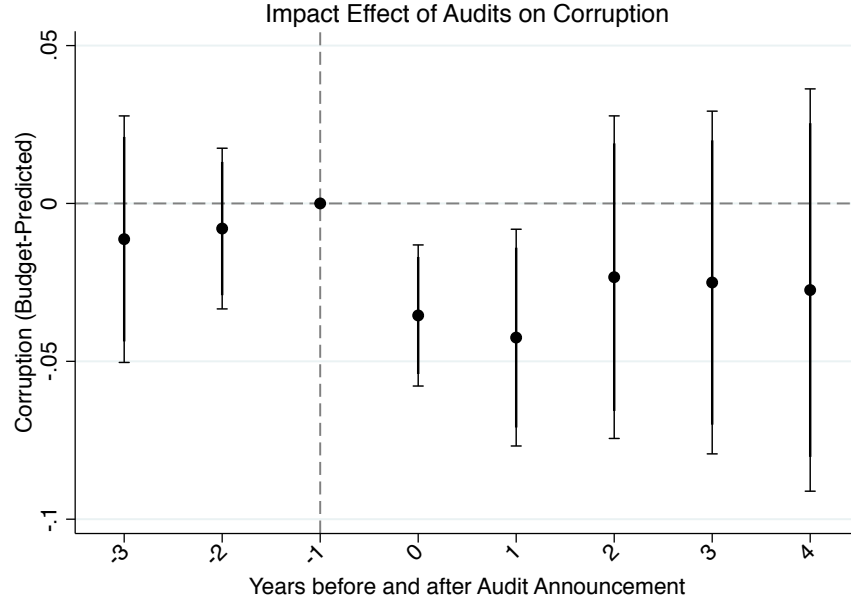
---

<sup>13</sup>Bobonis et al. (2016) study a similar research question in Puerto Rican municipalities. The authors focus on (non random) audit of municipal accounts, finding that audits do not persistently reduce corruption in that case.

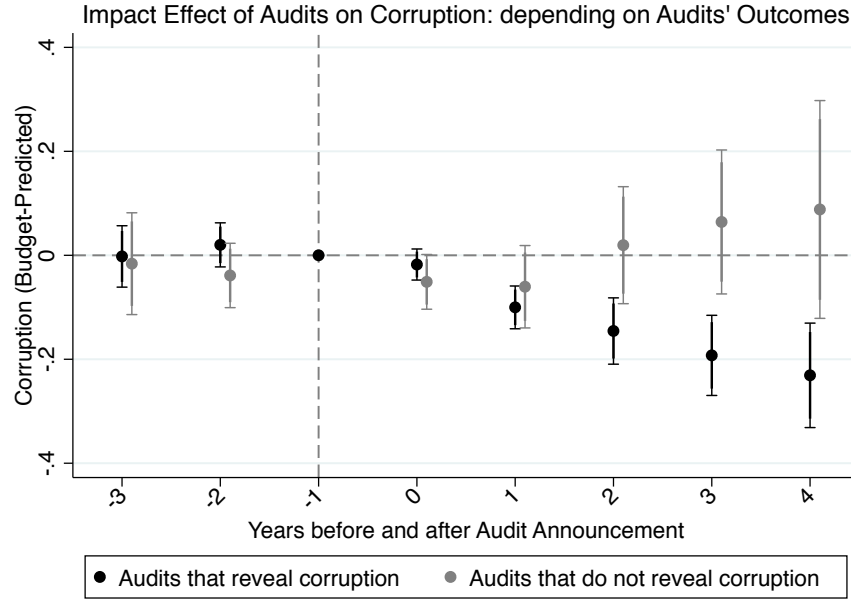
<sup>14</sup>The former group includes those cities in which the audits discovered a positive amount of corruption (measured with the variable narrow corruption), while the latter group includes those municipalities in which the audit did not find any type of corruption.

Figure 3: Dynamic Effect of Audits on Fiscal Corruption

(a) All municipalities



(b) Corrupted vs. Non-corrupted



*Notes:* Event study estimates for dynamic effect of audits on budget-predicted corruption. Error spikes give 95% confidence intervals, with standard error clustered by state. Top panel: all audits; bottom panel: audits that found corruption (in black); audits that did not find corruption (in grey). For the analysis on all audits leads are jointly insignificant (p-value=0.908) and lags are jointly significant (p-value=0.003). For the analysis on audits that found corruption leads are jointly insignificant (p-value=0.151) and lags are jointly significant (p-value=0.0003).

reduction in corruption, yet we cannot rule out that it consists of a substitution toward corrupted activities that are not detected by our model.

We report supporting analyses in Appendix C. First, we test the political accountability channel by checking whether the effect of the audit is stronger when local political competition is high. We find that the answer is yes: In cities where the mayor has been elected with a small margin of victory, the impact of the audit is stronger. Second, we check whether the main results may be explained by post-audit budget adjustments that might mechanically take place when a municipality is found to be corrupted. We show that this is not the case: The main results do not change if we control for total spending (per-capita) in the main regression specification, and we show that the occurrence of an audit does not affect future levels of municipal expenditures (per-capita).

Additional insights like these were not possible to obtain with the standard methods used in the previous literature. The number of multiple audits is too small, and the cross-sectional data too sparse, to analyze the rich comparative dynamics that we can do with our ML-predicted panel data. As with the analysis of revenue shocks, this second application on the audit effect on corruption demonstrates the usefulness of our machine learning approach to measuring corruption. The expanded datasets produced by machine prediction could be broadly useful for social scientists interested in corruption and governance.

## 5. Using Machine Learning to Guide Audit Policy

Besides extending datasets for empirical analysis, our machine predictions for corruption risk can also be used to guide policymakers. This section outlines a policy simulation for how corruption policy could be supported. We start with a baseline targeting policy based on predicted corruption risk, showing that targeted audits can detect more corruption than random audits. Second, we consider the issue of political bias in the risk scoring algorithm towards different mayor party affiliations, and analyze the performance of a politically neutral targeting policy. Third, we discuss additional caveats and complications with implementing a targeted audit system.

### 5.1. Targeted Audit Policy

To set the stage for targeted audits, let's first consider the performance of the status quo random audit system. Recall that there are 5563 municipalities in the dataset. In

our data, the agency audited 203 municipalities per year on average. We take that as the baseline for our policy analysis.

A set of key statistics from the random audits baseline are reported in Table 5 Column 1. We know from the audit data that the true corruption rate is 0.4664.<sup>15</sup> Since the audits are random, the corruption rate conditional on audit is still 0.4664. Taking the baseline of 203 audits per year, the audit probability (and therefore detection rate) is roughly 0.0365.<sup>16</sup> These numbers translate to about 95 corrupt municipalities detected for the average year of audits.

How can our machine predictions improve this outcome? We start by ranking the municipalities by corruption risk. That is, we apply the baseline gradient boosting model to the budget data for each municipality  $i$  from year  $t$  to produce  $\hat{y}_{it}$ . Then for each year  $t$ , we have an ordinal ranking of the municipalities (1 through 5563) by predicted probability of corruption. The proposed policy is to replace random audits with audits targeted by predicted corruption risk. Rather than sampling 203 municipalities uniformly from the distribution, the agency could audit the top 203 with the highest  $\hat{y}_{it}$ . These are municipalities that have a level of corruption probability higher than 0.847 in the average year.

This policy is illustrated in Figure 4. The horizontal axis gives the predicted corruption risk, and the blue marks give the true corruption rate at that risk level using the audit outcomes. The horizontal red dashed line at 0.037 gives the audit probability under random audits. The vertical green dashed line indicates the average threshold corruption risk (0.868, s.e. = .004) above which municipalities are targeted for audit.<sup>17</sup> The histogram indicates the distribution of the corruption risk predictions, with the top two bins containing the approximately 203 municipalities to be targeted.

Table 5 Column 2 reports statistics on the expected outcomes of this policy based on the true audit results.<sup>18</sup> Because the audited municipalities are higher risk, the

---

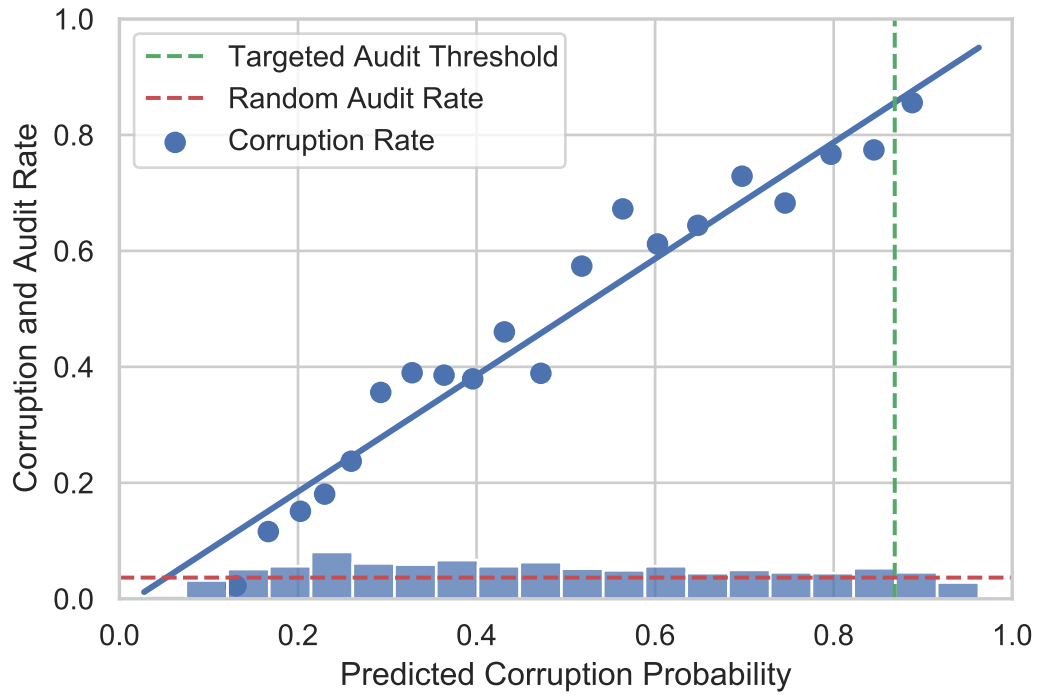
<sup>15</sup>Recall that the base rate in Section 3 was .422. The difference here is that we simplify the dataset to have a single observation per audit. In Section 3, the dataset included all fiscal years checked by the audit, which slightly changes the mean corruption rate.

<sup>16</sup>This is the unconditional probability. Because audits were randomly assigned within state, it could vary slightly across states. Using the unconditional probability simplifies the exposition without changing the qualitative implications.

<sup>17</sup>The threshold is an average across years, since the level of predicted corruption can vary across years.

<sup>18</sup>The statistics from Column 3 come from a politically neutral audit targeting policy. We revisit these numbers in the next subsection.

Figure 4: Targeted Auditing Based on Corruption Risk



*Notes:* Illustration of targeted auditing policy. Blue marks give the observed corruption rate from the audits (vertical axis) separately in 20 quantile bins constructed from the XGBoost model's predicted corruption probability (horizontal axis), with the blue diagonal giving line of best fit. Blue histogram shows the density of the predicted corruption probability. Red horizontal dashed line shows the audit probability under random audits ( $=0.037$ ). Vertical dashed green line shows the average threshold ( $=0.868$ ) above which a municipality is audited based on the targeting rule.

Table 5: Performance Metrics for Targeted Auditing Policies

	Random Audits (1)	Targeted Audits (2)	Fair Targeting (3)
Corruption Rate, if Audited	0.4664	0.8563 (0.0163)	0.8364 (0.0173)
Audit Rate, if Corrupt	0.0365	0.0671 (0.0013)	0.0655 (0.0014)
$\hookrightarrow$ Ratio over Random Audits		1.836 (0.035)	1.793 (0.037)

*Notes:* Metrics for comparing the effectiveness of audit policies: random audits (column 1), targeting audits to the municipalities with the highest corruption risk (column 2), or targeting audits with highest corruption with the constraint that all political parties are audited at the same rate. "Political party" means the set of municipalities where that party controls the mayor's office and includes PT, PMDB, PSDB, PTB, and DEM (formerly PFL). "Corruption Rate, if Audited" is the share of audited municipalities where narrow corruption is detected, for the respective policy. "Audit Rate, if Corrupt" is the expected probability of being audited, if narrow corrupt, under the various policies. Column 1 reports the observed rates in the data. In Columns 2 and 3, statistics give the mean and standard error (in parentheses) across five values for the predicted corruption risk, produced using different training-set folds. "Ratio over Random Audits" is the "Audit Rate, if Corrupt" value for the indicated policy, divided by that value under random audits.

audits are more effective: conditional on audit, the detected corruption rate of 0.856 is almost double ( $1.84\times$ ) that of the status quo policy (0.466).<sup>19</sup> Out of 203 audits, that corresponds to 168.2 corrupt municipalities detected, rather than 98.8.

Next, we consider the audit rate conditional on being corrupt. This value can be understood as the expected strength of enforcement or deterrence level. The conditional audit rate under targeting is 0.067, again almost 2x the status quo rate of 0.037. That is, corrupt mayors have a 6.6% chance of being discovered, rather than a 3.7% chance. Overall, targeting makes a big difference in policy effectiveness.

The numbers for targeted auditing indicate a significant policy improvement. For the same number of implemented audits (and presumably the same allocation of government resources), the targeted approach detects 84% more corrupt municipalities. Because successful audits reduce corruption (see Section 4.2 above), the targeted policy would

<sup>19</sup>Note that a counterfactual policy with the opposite goal (minimizing corruption detection) could target the lowest-risk municipalities and realize a detection rate of just 0.03.



also reduce the frequency of corrupt activities in Brazil. To achieve the same number of corruption detections as the status quo policy (95 municipalities), only 111 targeted audits are needed, down from 203 random audits. This decrease of 45%, or 91 audits per year, could imply a significant reduction in administrative expenditures.

To check robustness of these results, Appendix Table A9 reports analogous statistics to Table 5 for alternative specifications. First, we got statistically identical policy improvements using the model with alternative train/test splitting based on municipality rather than municipality-year. Second, analyzing a policy based on the alternative measure of corruption from [Avis et al. \(2018\)](#) obtained proportionally larger improvements on the status quo in terms of detecting corruption. Overall, the machine learning approach to support anti-corruption policy is robust to such implementation choices.

### 5.2. *Adjusting for Political Bias in Targeted Corruption Audits*

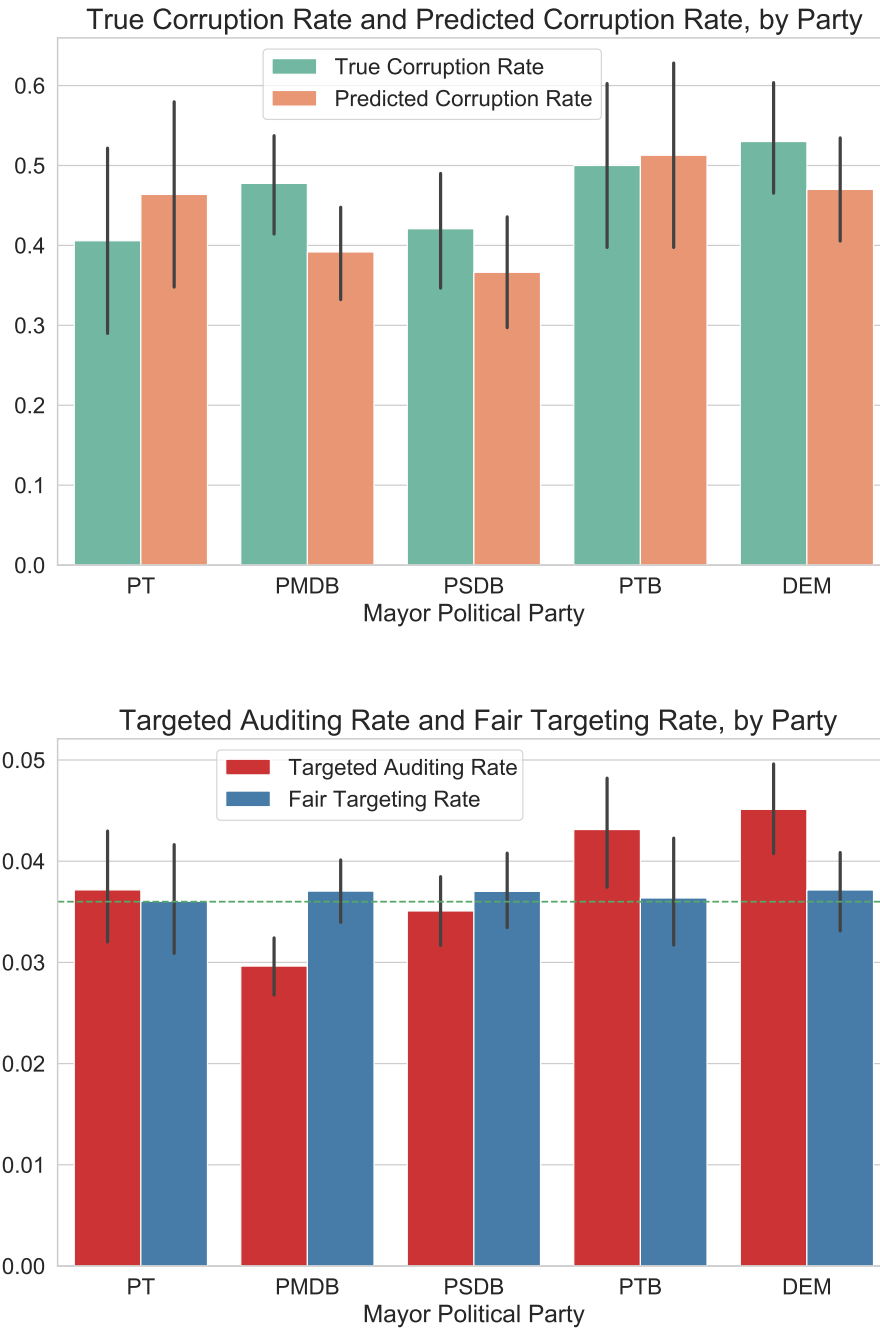
A key strength of randomized audits is that they are fair to all groups. A potential institutional barrier to the implementation of an AI-targeted anti-corruption system is the concern that it would be biased toward some political parties. In this section we consider whether the algorithm is biased toward different political parties, and if so, how to adjust the algorithm to produce fair outcomes.

We start by exploring how the corruption risk prediction algorithm treats the different political parties in our dataset. We focus on the five largest parties, as indicated by the average share of municipalities they control in our period. These parties, ranked roughly from most left-wing to most right-wing ([Power and Rodrigues-Silveira, 2019](#)), are PT, PMDB, PSDB, PTB, and DEM (formerly PFL). The distribution of municipality-terms by party is shown in Appendix Figure A2.

For each party, we compute the true corruption rate (from the random audits) and the predicted corruption rate (from the algorithm). These statistics are visualized in Figure 5, Top Panel. We can see that true corruption (green bars) varies somewhat across parties. For example, PSDB has a relatively low corruption rate, while PTB has a relatively high corruption rate. Although the differences across parties are mostly reproduced in the model’s predictions (orange bars), there is some important variation. For PMDB, in particular, the algorithm somewhat understates the risk of corruption relative to the true rate.

After seeing these numbers, politicians and policymakers may be skeptical about introducing targeted audits. What can be done to address this skepticism? Questions like

Figure 5: Corruption Risk and Targeted Auditing, by Party



*Notes:* Top Panel reports the true corruption rate in the audit data (in green bars) next to the predicted corruption rate from our XGBoost classifier (in orange bars), separately by the five political parties (meaning control of the mayor's office). Bottom Panel compares the auditing rates by party, under unconstrained targeting (red bars) and constrained targeting that equalizes audit rates across parties (blue bars). Horizontal dashed line gives the average audit rate in the sample. In both plots, 95% confidence interval spikes constructed by bootstrapping.

these are addressed by the burgeoning literature in algorithmic fairness (Barocas et al., 2019). This research subfield has promulgated formal definitions of fairness in automated decision-making and formulas for measuring them from datasets of predictions or decisions. The classic case study is criminal risk scoring, where existing algorithms have been shown to be biased toward racial minorities along some (but not all) definitions of fairness (Chouldechova, 2017; Berk et al., 2018; Kasy and Abebe, 2020).

This scholarship has developed a family of approaches for adjusting algorithmic decision procedures in order to mitigate fairness violations. An intuitive approach, which we follow here, is to separate the problem into a prediction step and a decision step. Rambachan et al. (2020) show that any and all equity concerns can be addressed solely at the decision stage, with the prediction stage being untouched. This *post-processing* approach is distinct from the more technically complex *pre-processing* or *constrained optimization* approaches that are explored in the computer science literature (see Barocas et al., 2019, ch. 3). The advantage of the latter methods is that the model does not need access to the sensitive covariate – normally, race/ethnicity – in order to produce a fair decision. In our setting, the sensitive covariate (city mayor party affiliation) is not that sensitive after all, and it will always be available in practice. Thus we take the post-processing approach.

Formally, we propose the following politically neutral targeting policy. As noted, the prediction algorithm is not changed at all. We start with  $\hat{y}_{it}$  for each municipality-year and the resulting corruption-risk ranking for all municipalities in a given year. Instead of taking the highest-ranked municipalities from the whole set, however, we produce separate rankings for each party. Within each party, we audit the same share of municipalities. Then by construction, the incidence of audits is equal across parties.

Figure 5 Bottom Panel shows the impact of fair targeting (blue bars) relative to unconstrained targeting (red bars). As intended, the fair audits have identical frequencies for each party (up to a rounding error). Comparing to the unconstrained rates, however, this fairness adjustment has significant redistributive consequences. On the one hand, PTB and DEM benefit from the introduction of fair targeting and are audited less often. On the other hand, fair targeting increases the audit risk for PMDB-controlled municipalities.

A second question is how fair targeting changes the overall effectiveness of audits, relative to unconstrained targeting. Revisiting Table 5, we see in Column 3 that the discovered corruption rate for audited municipalities is 0.836, still far higher than the

random baseline (0.4664). Discovered corruption is just slightly less in magnitude than the main targeting policy, and the difference is not statistically significant. In terms of deterrence – the audit rate, conditional on corruption – the nonpartisan policy still maintains significant policy effectiveness gains: 0.065, which is still  $1.793\times$  higher than the audit rate of 0.0365 under random assignment. This is quite close to, and not statistically different from, the unconstrained targeting policy. Overall, adjusting targeted anti-corruption policies to equalize audit rates across political parties does not significantly undermine the effectiveness of those policies.

### 5.3. *Additional Issues and Caveats*

Besides the incidence across political parties, there are a number of practical issues with a targeted auditing policy that would have to be addressed. First, the policy simulation considered so far has a single round of targeted audits. At least in the short run, multiple targeted audit rounds would be possible and effective if they used the public finances data from before the first audit. Subsequent to the first round of audits, however, the budget accounts would likely contain less information about corruption due to behavioral responses by local officials. The existing model, when applied to post-targeting accounting data, would likely produce significant errors that would favor the more savvy mayors. Still, it could reduce the net marginal benefit of corrupt activities by increasing the expected cost of corrupt fiscal actions that are not easily substitutable.<sup>20</sup>

In light of the behavioral responses, a question arises about how much information to publicize about audit targeting. One option would be to give full information about the policy and the associated model weights. This option would increase deterrence against corruption actions that are not easily substitutable. But it would reduce deterrence against substitutable actions, which could be easily gamed. Another option would be to start targeting audits without giving any information about how targeting is done. Presumably, over time corrupt officials could learn how municipalities are targeted, but it is unclear whether this could be done quickly enough to allow manipulation of accounts to avoid audits.

Understanding the relevance of these factors would require additional empirical evidence, preferably through randomized interventions. The specific numbers from our simulation should be taken with some skepticism, given the previous work showing that the

---

<sup>20</sup>Our setting is not amenable to the "manipulation-proof machine learning" method from [Björkegren et al. \(2020\)](#), which requires information on the cost function over corruption activities.

introduction of algorithms into decision-making can have smaller-than-expected effects (e.g. [Stevenson and Doleac, 2019](#)). There could be many reasons that an audit-targeting policy would not be as effective as outlined here.

In any case, a longer-term system of targeted audits would be more effective if some random audits were maintained. In such a mixed system, targeted audits would be used to detect and deter corruption for the highest-risk municipalities. Random audits would be maintained for two reasons. First, even apparently low-risk municipalities (including those who are good at fooling the algorithm) would have some chance of being audited and therefore face some deterrence incentive. Second, the results of the random audits would be used to update the algorithm parameters for guiding the next round of targeted audits. Determining the optimal mix of targeted and random audits would require more information and more assumptions on the deterrence effect of both types of audits.

## 6. Conclusion

This paper has shown that corruption in local governments can be reliably detected, predicted, and measured using public budget accounts data. We have shown that the resulting synthetic measurements can then be used in downstream empirical analysis, as we can produce the same empirical results using corruption predictions in municipalities that were never audited. Beyond expanding on empirical work, the corruption predictions can be used to guide policy responses to corruption. Our counterfactual policy estimates indicate substantial gains from such a policy, even when constraining the algorithm to treat each political party equitably.

This research adds to the emerging literature using machine learning and other tools from data science to explore new datasets and questions ([Kleinberg et al., 2015](#); [Athey, 2018](#)). Our method of detecting corruption has the potential to substantially expand the stock of datasets available for economists studying development, political economy, and public finance. Within Brazil, researchers will no longer be constrained to the relatively small set of municipalities that were audited. Outside of Brazil, the method could in principle be applied in any context with ground-truth labels for corruption. Something that can and should be explored is whether the corruption predictions produced in Brazil could be valid for other countries and settings.

## References

- Andini, M., E. Ciani, G. de Blasio, A. D’Ignazio, and V. Salvestrini (2018). Targeting with machine learning: An application to a tax rebate program in Italy. *Journal of Economic Behavior & Organization* 156, 86–102.
- Ash, E., M. Morelli, and R. Van Weelden (2017). Elections and divisiveness: Theory and evidence. *The Journal of Politics* 79(4), 1268–1285.
- Assunção, J., R. McMillan, J. Murphy, and E. Souza-Rodrigues (2019). Optimal environmental targeting in the Amazon rainforest. Technical report, National Bureau of Economic Research.
- Athey, S. (2018). The impact of machine learning on economics. In *The economics of artificial intelligence: An agenda*. University of Chicago Press.
- Athey, S. and S. Wager (2021). Policy learning with observational data. *Econometrica* 89(1), 133–161.
- Avis, E., C. Ferraz, and F. Finan (2018). Do government audits reduce corruption? estimating the impacts of exposing corrupt politicians. *Journal of Political Economy* 126(5), 1912–1964.
- Bandiera, O., A. Prat, S. Hansen, and R. Sadun (2020). CEO behavior and firm performance. *Journal of Political Economy* 0(0), 000–000.
- Bansak, K., J. Ferwerda, J. Hainmueller, A. Dillon, D. Hangartner, D. Lawrence, and J. Weinstein (2018). Improving refugee integration through data-driven algorithmic assignment. *Science* 359(6373), 325–329.
- Barocas, S., M. Hardt, and A. Narayanan (2019). *Fairness and machine learning: Limitations and Opportunities*.
- Battiston, P., S. Gamba, and A. Santoro (2020). Optimizing tax administration policies with machine learning. *University of Milan Bicocca Department of Economics, Management and Statistics Working Paper* (436).
- Becker, G. S. (1968). Crime and punishment: An economic approach. *Journal of Political Economy* 76(2), 169–217.
- Belloni, A., V. Chernozhukov, and C. Hansen (2014). High-dimensional methods and inference on structural and treatment effects. *Journal of Economic Perspectives* 28(2), 29–50.
- Berk, R., H. Heidari, S. Jabbari, M. Kearns, and A. Roth (2018). Fairness in crimi-

- nal justice risk assessments: The state of the art. *Sociological Methods & Research*, 0049124118782533.
- Björkegren, D., J. E. Blumenstock, and S. Knight (2020). Manipulation-proof machine learning. *arXiv preprint arXiv:2004.03865*.
- Bobonis, G. J., L. R. Cámara Fuertes, and R. Schwabe (2016). Monitoring corruptible politicians. *American Economic Review* 106(8), 2371–2405.
- Brollo, F., T. Nannicini, R. Perotti, and G. Tabellini (2013). The political resource curse. *American Economic Review* 103(5), 1759–96.
- Cavalcanti, F., G. Daniele, and S. Galletta (2018). Popularity shocks and political selection. *Journal of Public Economics* 165, 201–216.
- Chen, T. and C. Guestrin (2016). Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794.
- Cheol, L. and J. Mikesell (2018). The impact of public officials’ corruption on the size and allocation of u.s. state spending. *Public Administration Review*, 346–359.
- Chong, A., A. L. De La O, D. Karlan, and L. Wantchekon (2015). Does corruption information inspire the fight or quash the hope? a field experiment in mexico on voter turnout, choice, and party identification. *The Journal of Politics* 77(1), 55–71.
- Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data* 5(2), 153–163.
- Colonnelli, E., J. A. Gallego, and M. Prem (2019). What predicts corruption? *Available at SSRN 3330651*.
- Conley, T. G. and F. Decarolis (2016). Detecting bidders groups in collusive auctions. *American Economic Journal: Microeconomics* 8(2), 1–38.
- Coviello, D. and S. Gagliarducci (2017). Tenure in office and public procurement. *American Economic Journal: Economic Policy* 9(3), 59–105.
- Daniele, G. and T. Giommoni (2020). Corruption under austerity. *BAFFI CAREFIN Centre Research Paper No. 2020-131*.
- De Angelis, I., G. de Blasio, and L. Rizzica (2020). Lost in corruption. evidence from eu funding to southern italy. *Italian Economic Journal*, 1–23.
- Djankov, S., R. La Porta, F. Lopez-de Silanes, and A. Shleifer (2003). Courts. *The Quarterly Journal of Economics* 118(2), 453–517.
- Ferraz, C. and F. Finan (2008). Exposing corrupt politicians: The effects of Brazil’s

- publicly released audits on electoral outcomes. *The Quarterly Journal of Economics* 123(2), 703–745.
- Feurer, M., K. Eggensperger, S. Falkner, M. Lindauer, and F. Hutter (2018). Practical automated machine learning for the automl challenge 2018. In *International Workshop on Automatic Machine Learning at ICML*, pp. 1189–1232.
- Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, 1189–1232.
- Gallego, J., G. Rivero, J. D. Martínez, et al. (2018). Preventing rather than punishing: An early warning model of malfeasance in public procurement. Technical report.
- Gentzkow, M. and J. M. Shapiro (2010). What drives media slant? evidence from us daily newspapers. *Econometrica* 78(1), 35–71.
- Gentzkow, M., J. M. Shapiro, and M. Taddy (2019). Measuring group differences in high-dimensional choices: Method and application to congressional speech. *Econometrica* 87(4), 1307–1340.
- Glaeser, E. L., A. Hillis, S. D. Kominers, and M. Luca (2016). Crowdsourcing city government: Using tournaments to improve inspection accuracy. *American Economic Review* 106(5), 114–18.
- Hansen, S., M. McMahon, and A. Prat (2018). Transparency and deliberation within the fomc: a computational linguistics approach. *The Quarterly Journal of Economics* 133(2), 801–870.
- Hastie, T., R. Tibshirani, and J. Friedman (2009). *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media.
- Hessami, Z. (2014). Political corruption, public procurement, and budget composition: Theory and evidence from oecd countries. *European Journal of Political Economy* 34(C), 372–389.
- Hitsch, G. J. and S. Misra (2018). Heterogeneous treatment effects and optimal targeting policy evaluation. *Available at SSRN 3111957*.
- Kang, J. S., P. Kuznetsova, M. Luca, and Y. Choi (2013). Where not to eat? improving public policy by predicting hygiene inspections using online reviews. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 1443–1448.
- Kasy, M. and R. Abebe (2020). Fairness, equality, and power in algorithmic decision making. In *ICML Workshop on Participatory Approaches to Machine Learning*.



- Kleinberg, J., H. Lakkaraju, J. Leskovec, J. Ludwig, and S. Mullainathan (2018). Human decisions and machine predictions. *The quarterly journal of economics* 133(1), 237–293.
- Kleinberg, J., J. Ludwig, S. Mullainathan, and Z. Obermeyer (2015). Prediction policy problems. *American Economic Review* 105(5), 491–95.
- Knaus, M. C., M. Lechner, and A. Strittmatter (2018). Machine learning estimation of heterogeneous causal effects: Empirical monte carlo evidence. *The Econometrics Journal*.
- Knittel, C. R. and S. Stolper (2019). Using machine learning to target treatment: The case of household energy use. Technical report, National Bureau of Economic Research.
- Kyriacou, A. P., L. Muinelo-Gallo, and O. Roca-Sagalés (2015). Construction corrupts: Empirical evidence from a panel of 42 countries. *Public Choice* 165(1), 123–145.
- Lagaras, S., J. Ponticelli, and M. Tsoutsoura (2017). Caught with the hand in the cookie jar: Firm growth and labor reallocation after exposure of corrupt practices.
- Liu, C. and J. L. Mikesell (2019). Corruption and tax structure in american states. *The American Review of Public Administration* 49(5), 585–600.
- Liu, C., T. T. Moldogaziev, and J. L. Mikesell (2017). Corruption and state and local government debt expansion. *Public Administration Review* 77(5), 681–690.
- López-Iturriaga, F. J. and I. P. Sanz (2018). Predicting public corruption with neural networks: An analysis of spanish provinces. *Social Indicators Research* 140(3), 975–998.
- Machoski, E. and J. M. de Araujo (2020). Corruption in public health and its effects on the economic growth of brazilian municipalities. *The European Journal of Health Economics*, 1–19.
- Mauro, P. (1998). Corruption and the composition of government expenditure. *Journal of Public economics* 69(2), 263–279.
- Morris, S. D. and J. L. Klesner (2010). Corruption and trust: Theoretical considerations and evidence from mexico. *Comparative Political Studies* 43(10), 1258–1285.
- Mullainathan, S. and Z. Obermeyer (2019). A machine learning approach to low-value health care: wasted tests, missed heart attacks and mis-predictions. Technical report, National Bureau of Economic Research.
- Mullainathan, S. and J. Spiess (2017). Machine learning: an applied econometric approach. *Journal of Economic Perspectives* 31(2), 87–106.

- Olken, B. A. (2007). Monitoring corruption: evidence from a field experiment in indonesia. *Journal of political Economy* 115(2), 200–249.
- Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. (2011). Scikit-learn: Machine learning in python. *the Journal of machine Learning research* 12, 2825–2830.
- Power, T. J. and R. Rodrigues-Silveira (2019). Mapping ideological preferences in brazilian elections, 1994-2018: a municipal-level study. *Brazilian Political Science Review* 13(1).
- Rambachan, A., J. Kleinberg, S. Mullainathan, and J. Ludwig (2020). An economic approach to regulating algorithms. Technical report, National Bureau of Economic Research.
- Rockoff, J. E., B. A. Jacob, T. J. Kane, and D. O. Staiger (2011). Can you recognize an effective teacher when you recruit one? *Education finance and Policy* 6(1), 43–74.
- Stevenson, M. T. and J. L. Doleac (2019). Algorithmic risk assessment in the hands of humans. *Available at SSRN 3489440*.
- Widmer, P., S. Galletta, and E. Ash (2020). Media slant is contagious. *Center for Law & Economics Working Paper Series* 14.
- Winters, M. S. and R. Weitz-Shapiro (2013). Lacking information or condoning corruption: When do voters support corrupt politicians? *Comparative Politics* 45(4), 418–436.
- Zamboni, Y. and S. Litschig (2018). Audit risk and rent extraction: Evidence from a randomized evaluation in brazil. *Journal of Development Economics* 134, 133 – 149.

# A Machine Learning Approach to Analyze and Support Anti-Corruption Policy

## APPENDIX

## A. Additional Tables and Figures

Table A1: Balance sheets components

Year	Number of categories				Total
	Active	Passive	Expenditure	Revenue	
2001	56	46	43	52	197
2002	56	46	101	90	293
2003	57	48	100	90	295
2004	59	49	295	146	549
2005	63	52	298	151	564
2006	63	52	301	155	571
2007	64	52	309	170	595
2008	64	52	310	170	596
2009	80	57	331	198	666
2010	88	69	334	219	710
2011	89	69	335	219	712
2012	89	69	334	219	711

*Notes:* Summary tabulations on the number of components of the municipal budget by year and by macro category. The number of categories increases over time.

Table A2: Hyperparameter selection

Fold	L1 Penalty	L2 Penalty	Max Tree Depth	Learning Rate	Min. Child Weight	Tree Count
1	1	0.1	10	0.1	5	72
2	1	0.1	10	0.1	3	71
3	0.5	0.5	10	0.1	1	46
4	2	2	10	0.1	5	97
5	1	0.5	10	0.1	3	70

*Notes:* This table reports the hyperparameters selected for each of the 5 folds model training. Rows give the folds. L1 and L2 Penalty are regularization terms on the splitting decision that encourage smaller trees. Max Tree Depth is the max number of splits before a terminal node. Learning rate is how quickly parameters are updated during training. Minimum Child Weight is another regularization term, corresponding to the minimum number of observations required at each node. The last column is the number of trees grown in the resulting forest.

Table A3: Confusion Matrices

<i>Panel A. XGBoost</i>			
<i>Truth</i>	<i>Prediction</i>		
		Not Corrupt	Corrupt
	Not Corrupt	2573	485
	Corrupt	980	1261

---

<i>Panel B. OLS</i>			
<i>Truth</i>	<i>Prediction</i>		
		Not Corrupt	Corrupt
	Not Corrupt	1243	1815
	Corrupt	961	1280

---

<i>Panel C. LASSO</i>			
<i>Truth</i>	<i>Prediction</i>		
		Not Corrupt	Corrupt
	Not Corrupt	894	2164
	Corrupt	619	1622

---

<i>Panel D. Logistic regression</i>			
<i>Truth</i>	<i>Prediction</i>		
		Not Corrupt	Corrupt
	Not Corrupt	1568	1490
	Corrupt	840	1401

*Notes:* The table reports confusion matrices from the model predictions XGBoost (recall=0.562 and precision=0.722), OLS (recall=0.571 and precision=0.413), LASSO (recall=0.723 and precision=0.428) and Logistic regression (recall=0.625 and precision=0.484).

Table A4: Population thresholds for Inter-Government Transfers

Population interval	FPM coefficient
Below 10,189	0.6
10,189–13,584	0.8
13,585–16,980	1
16,981–23,772	1.2
23,773–30,564	1.4
30,565–37,356	1.6
37,357–44,148	1.8
44,149–50,940	2
Above 50,940	from 2.2 to 4

*Notes:* These coefficients have been introduced by *Decreto-lei* n. 1,881, 27 august 1981.

Table A5: Descriptive statistics for the Revenue Shocks Analysis

Population (1)	FPM transfers			N (5)
	Actual transfers (2)	Theoretical transfers (3)	Predicted Corruption (4)	
6,793 – 10,188	19.655	21.200	.442	1,429
10,189 – 13,584	25.642	28.771	.500	1,076
13,585 – 16,980	31.888	36.316	.527	805
16,981 – 23,772	38.445	44.019	.543	1,083
23,773 – 30,564	44.223	51.082	.529	629
30,565 – 37,356	50.869	58.113	.521	380
37,357 – 44,148	57.376	66.468	.510	253
44,149 – 50,940	62.389	72.368	.498	154
Total	33.440	37.930	.502	5,809

*Notes:* The sample includes all Brazilian municipalities with population in the interval 6,793-50,940. Population is the number of inhabitants. Actual and theoretical FPM transfers expressed in R\$100,000 at 2000 prices.

Table A6: Replication *Brollo et al. (2013)* with random samples

Random sample:	First (1)	Second (2)	Third (3)	Fourth (4)
<i>Panel A. First Stage</i>				
Theoretical transfers	0.7649*** (0.0215)	0.7100*** (0.0247)	0.6810*** (0.0485)	0.7344*** (0.0177)
<i>Panel B. Reduced Form</i>				
Theoretical transfers	0.0048*** (0.0007)	0.0042*** (0.0007)	0.0049*** (0.0007)	0.0038*** (0.0008)
<i>Panel C. 2SLS</i>				
Actual transfers	0.0063*** (0.0010)	0.0059*** (0.0010)	0.0072*** (0.0011)	0.0052*** (0.0010)
N. Observations	1115	1115	1115	1115

*Notes:* Effects of FPM transfers on (predicted) corruption measures. The four columns display the analysis focusing on four different random samples with 1,115 observations. Panel A reports the estimates of the first-stage analysis, the dependent variable is *actual transfers*. Panel B reports the estimates of reduced form analysis, the dependent variable is *predicted corruption*. Panel C reports the estimates of the 2sls estimates, the dependent variable is *predicted corruption* and *actual transfers* is instrumented with *theoretical transfers*. Column headings indicate the sample of municipalities included. All regressions controls for a third-order polynomial in normalized population size, term dummies, and macro-region dummies. Robust standard errors clustered at the municipal level are in parentheses: \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .



Table A7: Effect of Revenue Shocks on Corruption - Alternative Predictions

Dep. var.: Predicted corruption	All cities	
	Prediction demographics (1)	Prediction budget (without FPM) (2)
<i>Panel A. Reduced Form</i>		
Theoretical transfers	-0.0002 (0.0008)	0.0045*** (0.0003)
<i>Panel B. 2SLS</i>		
Actual transfers	-0.0003 (0.0011)	0.0065*** (0.0005)
N. Observations	5808	5808

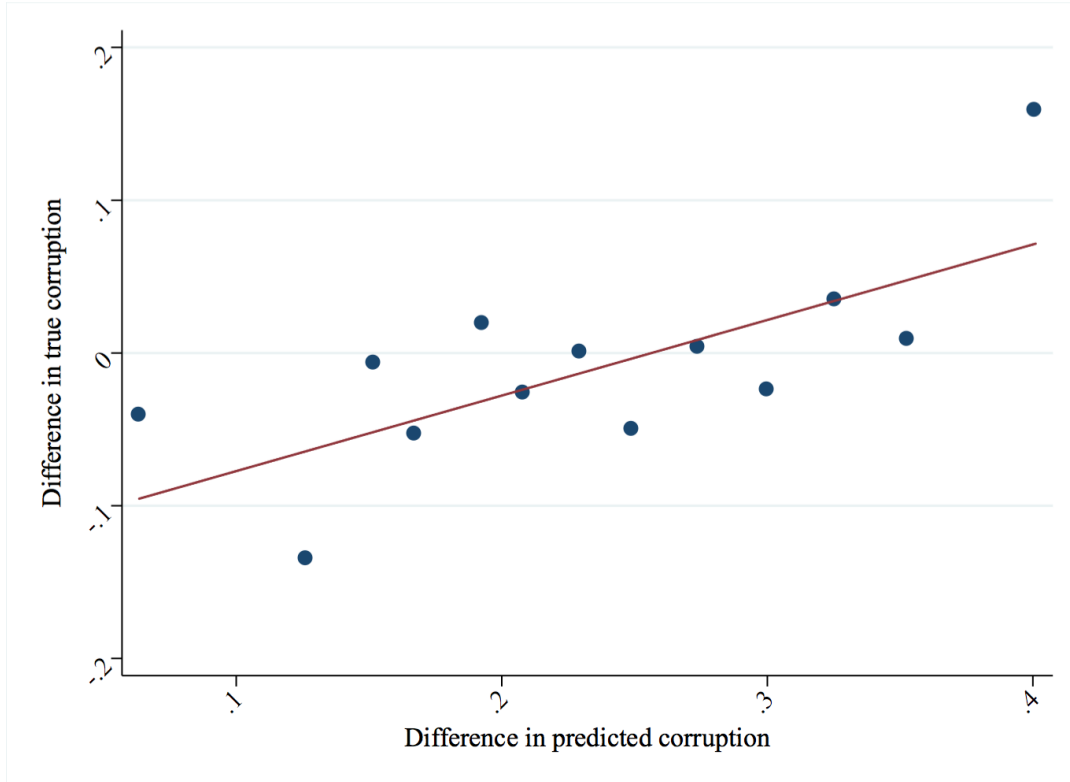
*Notes:* Effects of FPM transfers on (predicted) corruption measures: column (1) contains the analysis with the predictions built using as predictors a set of municipal demographic characteristics, and column (2) contains the analysis with the predictions built with budget predictors where FPM transfers are permuted randomly. Panel A reports the estimates of reduced form analysis, the dependent variable is *predicted corruption*. Panel B reports the estimates of the 2sls estimates, the dependent variable is *predicted corruption* and *actual transfers* is instrumented with *theoretical transfers*. The sample includes all Brazilian municipalities. All regressions controls for a third-order polynomial in normalized population size, term dummies, and macro-region dummies. Robust standard errors clustered at the municipal level are in parentheses: \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

Table A8: Coefficient Estimates for Event Study Analysis

	All cities (1)	Cities with corruption (2)	Cities without corruption (3)
Year pre4 and behind	-0.0171 (0.0245)	-0.0287 (0.0427)	-0.0052 (0.0748)
Year pre3	-0.0118 (0.0190)	-0.0024 (0.0287)	-0.0164 (0.0476)
Year pre2	-0.0078 (0.0124)	0.0203 (0.0205)	-0.0390 (0.0302)
Audit year	-0.0358*** (0.0109)	-0.0177 (0.0145)	-0.0506* (0.0254)
Year post1	-0.0429** (0.0166)	-0.1002*** (0.0200)	-0.0597 (0.0387)
Year post2	-0.0238 (0.0246)	-0.1456*** (0.0311)	0.0205 (0.0545)
Year post3	-0.0253 (0.0262)	-0.1924*** (0.0376)	0.0659 (0.0672)
Year post4	-0.0276 (0.0308)	-0.2307*** (0.0490)	0.0903 (0.1018)
Year post5	-0.0156 (0.0418)	-0.2585*** (0.0620)	0.1581 (0.1185)
Years post6 and more	-0.0364 (0.0478)	-0.3260*** (0.0711)	0.1756 (0.1294)
N. Observations	17252	8895	3086
Adjusted $R^2$	0.535	0.510	0.538

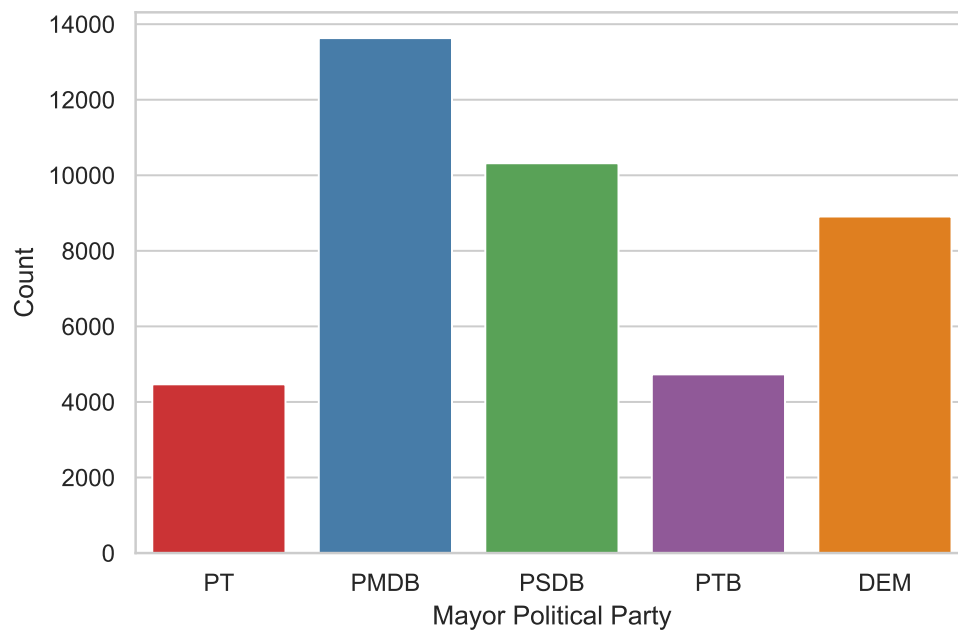
*Notes:* The dependent variable is (predicted) corruption measure - binary. The sample includes all the cities that receive an audit for the period 2001-2012. Column (1) includes the complete sample, Column (2) includes the sample of cities in which the audit discovered corruption (according to the definition of narrow corruption) and Column (3) includes the sample of cities in which the audit did not discover any type of corruption. The specification includes city and year fixed effects. Robust standard errors clustered at the state level are in parentheses: \*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$ .

Figure A1: Difference in true and predicted corruption for cities audited twice



*Notes:* The figure focuses on cities that have been audited twice and it shows a binscatter between the difference over time in the true levels of corruption using the data from [Brollo et al. \(2013\)](#) and the predicted levels of corruption. The analysis includes the following list of fixed effects and controls: first audit year and second audit year fixed effects, mean income, share of population employed, sector of occupation (agriculture, industry, commerce, transportation, services and public administration), share with college education, poverty rate, and Gini Coefficient of income. The coefficient of the corresponding regression is 0.495 (p-value 0.095).

Figure A2: Distribution of Party Control of Municipalities



*Notes:* Number of municipality-year observations for each party, in terms of the affiliation of the mayor in that municipality.

Table A9: Performance Metrics for Targeted Auditing Policies, Additional Specifications

	<i>Brollo et al (2013) Corruption Labels</i>			<i>Avis et al (2018) Corruption Labels</i>		
	Lottery	Targeted Audits		Lottery	Targeted Audits	
	(1)	(2)	(3)	(4)	(5)	(6)
Train/Test Split		Muni-Year	Muni		Muni-Year	Muni
Corruption Rate, if Audited	0.4664	0.8563 (0.0163)	0.8231 (0.0217 )	0.19	0.6743 (0.0186)	0.6725 (0.0257)
# Corrupt Munis Detected	94.8	174.0861 (3.3225)	167.3261 (4.45)	40.4	137.0670 (3.7793)	136.7004 (5.2294)
Audit Rate, if Corrupt	0.036	0.0671 (0.0013)	0.0645 (0.0017)	0.036	0.1241 (0.0034)	0.1237 (0.0047)
$\hookrightarrow$ Ratio to Lottery		1.836 (0.035)	1.7648 (0.0465)		3.3954 (0.0936)	3.3863 (0.1295)
Min Audit # Equivalent		110.9143 (2.0543)	115.5913 (3.0415)		60.1115 (1.7429)	60.4993 (2.4499)

*Notes:* Metrics for comparing the effectiveness of audit policies. Columns 1 through 3 use the main label of corruption from [Brollo et al. \(2013\)](#). Columns 4 through 6 use the alternative label of corruption from [Avis et al. \(2018\)](#). Columns 1 and 4 report the true rates under random audits. Columns 2, 3, 5, and 6 report the results from targeting audits, with Columns 2 and 5 using the main train/test sampling by municipality-year, and Columns 3 and 6 using the alternative grouped splitting by municipality. The rows report the different policy outcomes. "Corruption Rate, if Audited" is the share of audited municipalities where narrow corruption is detected. "# Corrupt Munis Detected" is the number of corrupt municipalities detected, out of the 203 audits implemented. "Audit Rate, if Corrupt" is the expected probability of being audited if corrupt. "Ratio to lottery" is the "Audit Rate, if Corrupt" value for the indicated policy, divided by that value under random audits. "Min Audit # Equivalent" is the number of audits needed under targeting to detect the same number of corrupt municipalities detected under the lottery system. For the audit-targeting statistics, we report the mean and standard error (in parentheses) across five values for the predicted corruption risk, produced using different training-set folds.

## B. Alternative Prediction Specifications

This appendix reports the performance metrics from some alternative corruption prediction specifications. First, to compare XGBoost model performance using not only budget factors but also fixed demographic factors, we apply random splits between training and test set by municipality, instead of by municipality-year. Appendix Table A10 shows the relative performance when we use budget data (column 1), when we add demographic characteristics (column 2), or when we use only demographic characteristics (column 3). The more conservative sampling specification in Column 1 reduces accuracy compared to the main-text specification, but it is still capturing significant predictive signal (in Column 3, AUC-ROC = 0.636 with budget and demographics). Comparing Column 1 to Column 2, we see that budget information is more predictive of corruption than demographic information.

Second, we replicate our predictive results by using the corruption measure from [Avis et al. \(2018\)](#). It is important to stress that there are structural differences between these two original measures of corruption. A first important difference is that the alternative measure is continuous, rather than binary. We have for each audited municipality the share of inspection orders that presented irregularities. The second difference is that we are missing the first audits, as we have information only from July 2006 through March 2013 (lotteries 22–38). Third, differently from the main main measure, with the alternative measure we do not know the exact year (or term) in which the irregularity took place. To overcome this limitation, we treat as audited the three years before the actual audit took place. Finally, to create a binary label from the continuous variable we identified as corrupted those municipalities with a share of irregularities in the top quartile of the distribution.

Despite these differences, Figure A3 shows that the predictions using the alternative corruption label are similar to those from the main analysis. They show similar rankings on average. The performance metrics are reported columns (4-7) of Appendix Table A10. Again, we find that XGBoost outperforms all the other methods. Indeed, we find accuracy metrics that are higher than those from the main analysis.

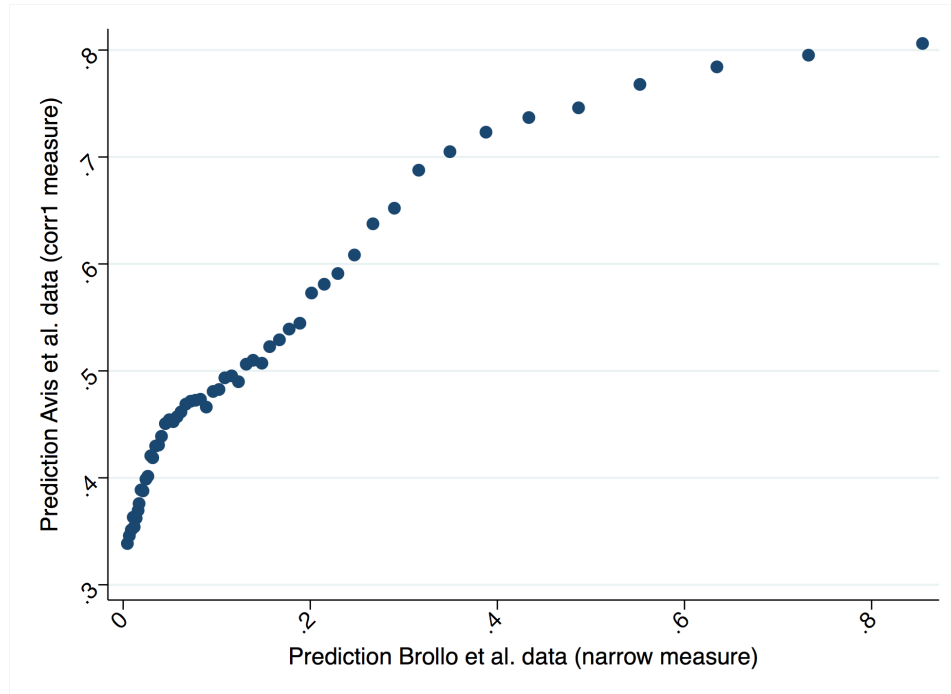
Finally, we find that most of our empirical results still holds when using the predictions from this alternative measure of corruption.

Table A10: Additional models performance

	XGBoost (municipal sampling)			Avis et al. (2018) data			
	Budget (1)	Budget + Demo (2)	Demo (3)	XGBoost (4)	OLS (5)	LASSO (6)	Logistic (7)
Accuracy	0.613 (0.012)	0.613 (0.004)	0.576 (0.009)	0.851 (0.005)	0.431 (0.048)	0.419 (0.062)	0.688 (0.036)
AUC-ROC	0.618 (0.015)	0.636 (0.006)	0.589 (0.012)	0.903 (0.009)	0.443 (0.065)	0.519 (0.033)	0.657 (0.010)
F1	0.486 (0.018)	0.498 (0.007)	0.476 (0.009)	0.635 (0.017)	0.311 (0.050)	0.375 (0.028)	0.485 (0.021)

*Notes:* The table provides the mean and standard error (in parentheses) across five values for the prediction performance, produced using different training-set folds. In columns (1-3) we use XGBoost models with municipal sampling, and different sets of predictors: only budget components in column (1), budget components and demographic characteristics in column (2) and only demographic characteristics in column (3). In columns (4-8) we report the predictions performance as in Table 2, but using the corruption data from Avis et al. (2018).

Figure A3: Predictions from Avis et al. (2018) vs. Predictions from Brollo et al. (2013)



*Notes:* The figure shows a biscatter between the predictions formed using the data from Avis et al. (2018) and the ones formed using the data from Brollo et al. (2013) for all municipality-year. The correlation between the two variables is 0.531.

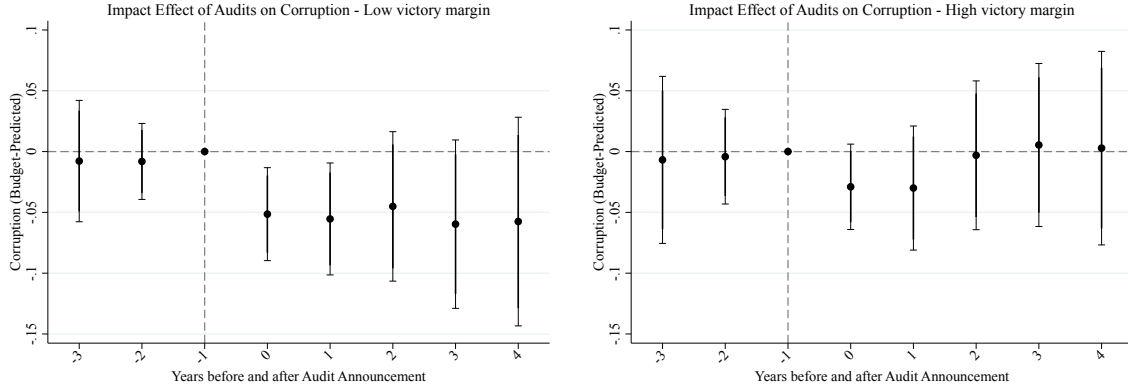
### C. Additional material for the effect of audits on corruption

In this Appendix we discuss a series of additional results for the event study analysis. First, we test the channel of political accountability. In particular, we aim to study whether the effect of the audit on future corruption is stronger where local political accountability is high and we focus on the variable margin of victory. This test is shown in Figure A4, which reports the analyses conducted with the full sample. The figures show that the effect is stronger in cities where the mayor won with a small margin of victory – below the median level – compared to cities where she won with a high margin – above the median level. This result suggests that the audit has a larger impact where the electoral competition is more pronounced. Overall, these results provide some evidence that political accountability affects the impact of an audit on future corruption.

Second, we check whether post-audit budget adjustments may explain the decline in predicted corruption levels after the audit. We provide two tests. First, we estimate the main model controlling for total expenditure, expressed in per-capita terms. This test is reported in Figure A5 and the results are similar to the ones of the main model, reported in figure 4. Secondly, we estimate the main model using as dependent variable the amount of total expenditure (per-capita): Figure A6 shows this test and it suggests that the audit does not have any significant effect on future levels of total spending. This result holds for the full sample and for the sample of corrupted and non-corrupted cities.

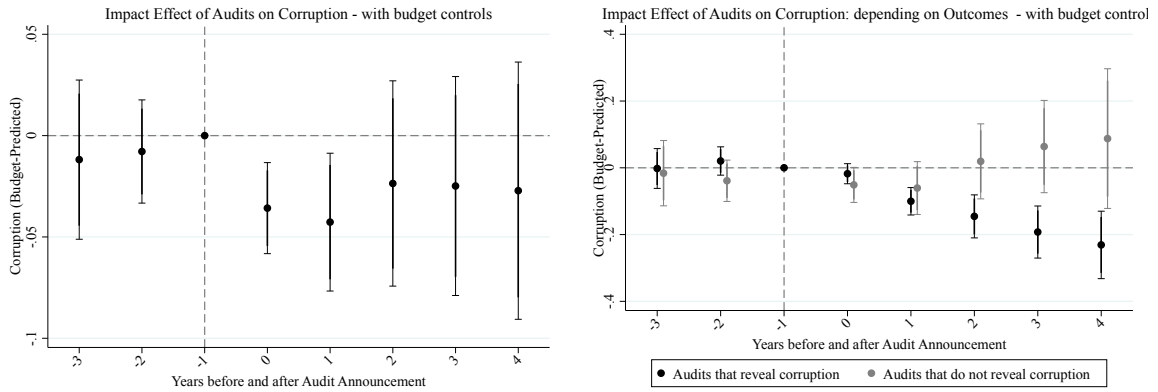


Figure A4: Dynamic effect of the audits - Margin of victory



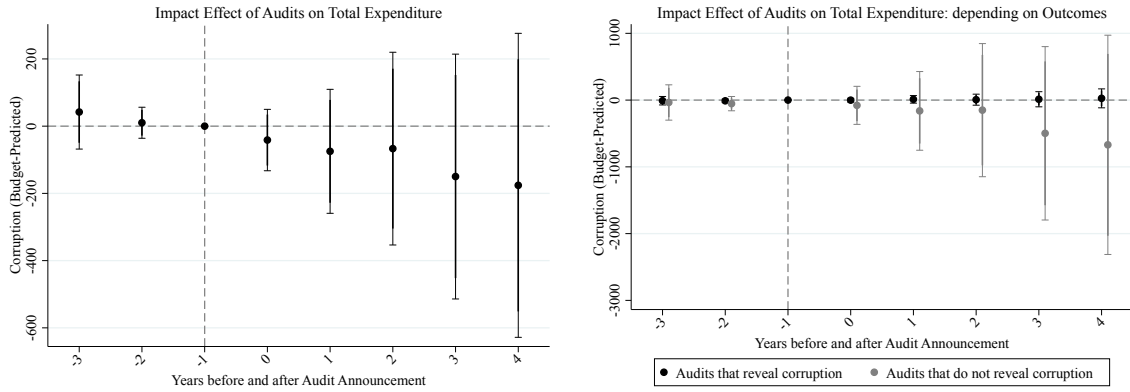
*Notes:* Event study estimates for dynamic effect of audits on budget-predicted corruption. Error spikes give 95% confidence intervals, with standard error clustered by state. In the left panel are considered only municipalities where the mayor won with a low margin of victory (below the median); In the right panel are considered only municipalities where the mayor won with a high margin of victory (above the median)

Figure A5: Dynamic effect of the audits - Controlling for total expenditure



*Notes:* Event study estimates for dynamic effect of audits on budget-predicted corruption. Error spikes give 95% confidence intervals, with standard error clustered by state. Left panel: all audits; right panel: audits that found corruption (in black); audits that did not find corruption (in grey). This regressions include as additional control municipal total expenditure.

Figure A6: Dynamic effect of the audits on total expenditure



*Notes:* Event study estimates for dynamic effect of audits on municipal total expenditure. Error spikes give 95% confidence intervals, with standard error clustered by state. Left panel: all audits; right panel: audits that found corruption (in black); audits that did not find corruption (in grey).