



THE LONDON SCHOOL
OF ECONOMICS AND
POLITICAL SCIENCE ■

Brunori, Paolo , Hufe, Paul & Mahler, Daniel (2023) The roots of inequality: estimating inequality of opportunity from regression trees and forests. *Scandinavian Journal of Economics*, 125(4), 900 - 932. <https://doi.org/10.1111/sjoe.12530>

<https://researchonline.lse.ac.uk/id/eprint/118220/>

Version: Published Version

Licence: [Creative Commons: Attribution 4.0](#)

[LSE Research Online](#) is the repository for research produced by the London School of Economics and Political Science. For more information, please refer to our [Policies](#) page or contact lseresearchonline@lse.ac.uk

Scand. J. of Economics 000(0), n/a, 2023

DOI: 10.1111/sjoe.12530

The roots of inequality: estimating inequality of opportunity from regression trees and forests*

Paolo Brunori[†]

London School of Economics, London, WC2A 2AE, UK

paolo.brunori@unifi.it

Paul Hufe[‡]

University of Bristol, Bristol, BS8 1TU, UK

paul.hufe@bristol.ac.uk

Daniel Mahler

World Bank, Washington, DC 20433, USA

dmahler@worldbank.org

Abstract

We propose the use of machine learning methods to estimate inequality of opportunity and to illustrate that regression trees and forests represent a substantial improvement over existing approaches: they reduce the risk of ad hoc model selection and trade off upward and downward bias in inequality of opportunity estimates. The advantages of regression trees and forests are illustrated by an empirical application for a cross-section of 31 European countries. We show that arbitrary model selection might lead to significant biases in inequality of opportunity estimates relative to our preferred method. These biases are reflected in both point estimates and country rankings.

Keywords: Equality of opportunity; machine learning; random forests

JEL classification: C38; D31; D63

*We thank Chiara Binelli, Marc Fleurbaey, Torsten Hothorn, Niels Johannesen, Andreas Peichl, Giuseppe Pignataro, Dominik Sachs, Jan Stuhler, Dirk Van de gaer, and Achim Zeileis for useful comments and suggestions. Furthermore, we are grateful for the comments received from seminar audiences at Princeton University, the University of Perugia, the University of Essex, the World Bank, ifo Munich, the University of Copenhagen, Canazei Winter School 2018, the European Commission JRC at Ispra, the EBE Meeting 2018, IIPF 2018, and the Equal Chances Conference in Bari. Any errors remain our own.

[†]Also affiliated with the University of Florence, Italy.

[‡]Also affiliated with IZA and CESifo.

© 2023 The Authors. *The Scandinavian Journal of Economics* published by John Wiley & Sons Ltd on behalf of Föreningen för utgivande av the *SJE*.

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

1. Introduction

Equality of opportunity is an important ideal of distributive justice. It has widespread support among the general public and its realization has been identified as an important goal of public policy intervention (Cappelen et al., 2007; Corak, 2013; Chetty et al., 2016; Alesina et al., 2018). In spite of its popularity, it is notoriously difficult to provide empirical estimates of equality of opportunity. Next to normative dissent about the precise factors that should be viewed as contributing to unequal opportunities, current estimation approaches are encumbered by ad hoc model selection that leads researchers to overestimate or underestimate inequality of opportunity.

In this paper, we propose the use of machine learning methods to overcome the issue of ad hoc model selection. Machine learning methods allow for flexible models of how unequal opportunities come about while imposing statistical discipline through criteria of out-of-sample replicability. These features serve to establish estimates of inequality of opportunity that are less prone to upward or downward bias.

The empirical literature on the measurement of unequal opportunities has been flourishing since the ground-breaking contribution by Roemer (1998), *Equality of Opportunity*. At the heart of Roemer's formulation is the idea that individual outcomes are determined by two sorts of factors: those factors over which individuals have control, which he calls "effort", and those factors for which individuals cannot be held responsible, which he calls "circumstances". While outcome differences due to effort exertion are morally permissible, differences due to circumstances are inequitable and call for compensation.¹ Grounded on this distinction, measures of inequality of opportunity quantify the extent to which individual outcomes are predicted by circumstance characteristics. They are usually calculated in a two-step procedure. First, researchers predict an outcome of interest from observable circumstances. Second, they calculate inequality in the distribution of predicted outcomes: the more predicted outcomes diverge, the more circumstances are associated with outcomes, and there is more inequality of opportunity.

Current approaches to estimate inequality of opportunity suffer from biases that are the consequence of critical choices in model selection. First, researchers have to decide which circumstance variables to consider for estimation.² The challenge of this task grows with the increasing availability

¹The distinction between circumstances and efforts underpins many prominent branches of the economics literature, such as the ones on intergenerational mobility (Chetty et al., 2014a,b), the gender pay gap (Blau and Kahn, 2017), and racial differences (Kreisman and Rangel, 2015). For different notions of equality of opportunity, see Arneson (2018).

²Roemer does not provide a fixed list of circumstance variables. Instead, he suggests that the set of circumstances should evolve from a political process (Roemer and Trannoy, 2015).

of high-quality datasets that provide very detailed information with respect to individual circumstances (Björklund et al., 2012; Hufe et al., 2017). On the one hand, discarding relevant circumstances from the estimation model limits the explanatory scope of circumstances and leads to downward-biased estimates of inequality of opportunity (Ferreira and Gignoux, 2011). On the other hand, including too many circumstances overfits the data and leads to upward-biased estimates of inequality of opportunity (Brunori et al., 2019). Second, researchers must choose a functional form according to which circumstances co-produce the outcome of interest. For example, it is a well-established finding that the influence of socio-economic disadvantages during childhood on life outcomes varies by biological sex (Dahl and Lochner, 2012; Chetty et al., 2016). In contrast to such evidence, many empirical applications presume that the effect of circumstances on individual outcomes is log-linear and additive while abstracting from possible interaction effects (Bourguignon et al., 2007; Ferreira and Gignoux, 2011). On the one hand, restrictive functional form assumptions limit the ability of circumstances to explain variation in the outcome of interest and thus force a downward bias on inequality of opportunity estimates. On the other hand, limitations in the available degrees of freedom might prove a statistically meaningful estimation of complex models with many parameters infeasible.

This discussion highlights the non-trivial challenge of selecting the appropriate model for estimating inequality of opportunity. Researchers must balance different sources of bias while avoiding ad hoc solutions. While this task is daunting for the individual researcher, it is a standard application for machine learning algorithms that are designed to make out-of-sample predictions of a dependent variable based on a number of observable predictors. In this paper, we use conditional inference regression trees and forests to estimate inequality of opportunity (Hothorn et al., 2006). Introduced by Morgan and Sonquist (1963) and later popularized by Breiman et al. (1984); Breiman (2001), they belong to a set of machine learning methods that is increasingly integrated into the statistical toolkit of economists (Varian, 2014; Mullainathan and Spiess, 2017; Athey, 2018). Trees and forests obtain predictions by drawing on a clear-cut algorithm that imposes only minimal assumptions about which circumstances interact in shaping individual opportunities, and how. Thereby, they restrict judgment calls of the researcher and inform model specification by data analysis. As a consequence, they cushion downward bias by flexibly accommodating different ways of how circumstance characteristics shape the distribution of outcomes. Moreover, the conditional inference algorithm branches trees (and constructs forests) by a

In empirical implementations, typical circumstances include biological sex, socio-economic background, and race.

4 Estimating inequality of opportunity from regression trees and forests

sequence of hypothesis tests that prevents the inclusion of noisy circumstance parameters. This feature reduces the potential for upward-biased estimates of inequality of opportunity through model overfitting. Hence, regression trees and forests address the detrimental consequences of ad hoc model selection in a way that is sensitive to both upward and downward bias in inequality of opportunity estimates.

To showcase the advantages of regression trees and forests, we compare them to existing estimation approaches in a cross-sectional dataset of 31 European countries. We demonstrate that current estimation approaches overfit (underfit) the data, which in turn leads to upward(downward)-biased estimates of inequality of opportunity. These biases are sizable. For example, some standard methods overestimate inequality of opportunity in the Nordic countries while they underestimate the extent of inequality of opportunity in Germany and France. As a consequence, these countries appear close in terms of their opportunity characteristics. Hence, standard estimation approaches can yield misleading information about the level of inequality of opportunity in different societies to policymakers and the general public alike.

While we demonstrate the advantages of regression trees and forests for estimations of inequality of opportunity, they are not a panacea to empirical challenges in this literature. First, regression trees and forests cannot address one of the most relevant sources of downward bias in inequality of opportunity estimates: missing data on relevant circumstances. Second, although regression trees and forests are less prone to upward and downward bias, the remaining bias can nevertheless be substantial when samples are small. Therefore, we encourage applied researchers to exercise caution when estimating inequality of opportunity on data with small number of observations.

The remainder of this paper is organized as follows. In Section 2, we give a brief introduction to current empirical approaches in the literature on inequality of opportunity. In Section 3, we introduce conditional inference regression trees and forests, and illustrate how to use them in the context of inequality of opportunity estimations. Section 4 contains an empirical illustration based on simulated data and the EU Survey of Income and Living Conditions. In this section, we also highlight the particular advantages of tree- and forest-based estimation methods by comparing them with the prevalent estimation approaches in the literature. We conclude in Section 5.³

³In a parallel paper, Blundell and Risa (2019) apply machine learning methods to the estimation of intergenerational mobility. In particular, they assess the completeness of rank–rank estimates of intergenerational mobility as measures of equal opportunities. In contrast to their work, we directly estimate inequality of opportunity statistics. Therefore, our focus is not on downward bias that follows from focusing on one circumstance only (i.e., parental income) but on balancing

2. Empirical approaches to equality of opportunity

2.1. Theoretical set-up and notation

Consider a population $\mathcal{N} = \{1, \dots, N\}$ and an associated vector of non-negative incomes $y = (y_1, \dots, y_N)$. Income is determined by two sets of factors: circumstances beyond individual control and individual efforts. We define the $(P \times 1)$ -vector $\omega_i \in \Omega$ as a comprehensive description of the circumstances of $i \in \mathcal{N}$. Analogously we define the $(Q \times 1)$ -vector $\theta_i \in \Theta$ as a comprehensive description of the efforts that are exerted by $i \in \mathcal{N}$. The income-generating function can be defined as

$$y = d(\omega, \theta). \quad (1)$$

Based on the realizations of individual circumstances, the population can be partitioned into types. We define the type partition $\mathcal{T} = \{t_1, \dots, t_M\}$, such that individuals are member of one type if they share the same circumstances: $i, j \in t_m \Leftrightarrow \omega_i = \omega_j$.

2.2. Measurement

Opportunity egalitarians are averse to inequalities that are rooted in circumstances; however, they are indifferent to inequalities that originate from individual effort exertion. In spite of the intuitive appeal of this idea, the literature has suggested a variety of formulations that differ in their precise normative content; see Ramos and Van de gaer (2016) for an overview. In this work, we exclusively focus on *ex ante* utilitarian measures of inequality of opportunity (Van de gaer, 1993; Checchi and Peragine, 2010). These are the most widely applied formulations in the empirical literature.⁴

According to the *ex ante* utilitarian view, the value of a type's opportunity set is pinned down by the expected value of its outcomes, $\mathbb{E}[y|\omega]$. Thus, the distribution of opportunities in a population can be expressed by the following counterfactual distribution y^C :

$$y^C = (y_1^C, \dots, y_i^C, \dots, y_N^C) = (\mathbb{E}[y_1|\omega_1], \dots, \mathbb{E}[y_i|\omega_i], \dots, \mathbb{E}[y_N|\omega_N]). \quad (2)$$

From this distribution, one can construct *ex ante* utilitarian measures of inequality of opportunity by choosing any functional $I(\cdot)$ that satisfies the following two properties:

both downward and upward bias if the set of available circumstances is large in relation to a given sample size.

⁴The use of machine learning methods is not restricted to *ex ante* utilitarian formulations and can be easily extended to alternative measures of inequality of opportunity.

6 Estimating inequality of opportunity from regression trees and forests

1. $I(y^C)$ decreases (increases) through transfers from i to j if i is from a circumstance type with a higher (lower) expected value of outcomes than the recipient j ;
2. $I(y^C)$ remains unaffected by transfers from i to j if they are members of the same type.

In most empirical applications, $I(\cdot)$ represents an inequality index satisfying the standard properties of anonymity, the principle of transfers, population replication, and scale invariance (Cowell, 2016).⁵ Examples of the latter are the Gini index or any member of the generalized entropy class. Note that the choice of $I(\cdot)$ is normative in itself as it specifies the extent of inequality aversion at different points of the counterfactual distribution y^C . For example, the mean logarithmic deviation (MLD) values compensating transfers to the most disadvantaged types more than the Gini index. In this work, we are agnostic about the normatively correct choice of $I(\cdot)$. While we present our main results in terms of the Gini index, we provide robustness checks based on other inequality indices in Section S.6 of the Supplementary Material.

2.3. Estimation

Given the measurement decisions described above, we require an estimate of the conditional distribution y^C . The data-generating process (DGP) described in equation (1) can be rewritten as

$$y = d(\omega, \theta) = f(\omega) + \epsilon = \mathbb{E}(y|\omega) + \epsilon. \quad (3)$$

Here, $\mathbb{E}(y|\omega)$ captures unfair variation due to observed circumstances. The independent and identically distributed error term ϵ captures both fair (individual effort) and unfair (unobserved circumstances) determinants of individual outcomes; hence, resulting measures of inequality of opportunity have a lower bound interpretation.

Estimating y^C is a prediction task in which the researcher tries to answer the following question: what outcome y_i do we expect for an individual who faces circumstances ω_i ? The precise form of $f(\cdot)$ is *a priori* unknown. In the vast majority of empirical applications, researchers address this lack of

⁵The β coefficient from intergenerational mobility regressions can also be interpreted as an *ex ante* utilitarian measure of inequality of opportunity. In the intergenerational mobility framework, $\beta = E(y_{ic}|y_{ip})/y_{ip}$, where y_{ip} represents parental income as the sole circumstance. Hence, the functional applied to the distribution of conditional expectations can be written as $I(\cdot) = 1/y_{ip}$. Note that β decreases (increases) through transfers from children from advantaged (disadvantaged) backgrounds to children from less (more) advantaged backgrounds. However, β remains unaffected by transfers between children from parental households with equal y_{ip} .

knowledge by invoking strong functional form assumptions. For example, they perform a log–linear regression of the outcome of interest on the set of observed circumstances and construct an estimate for y^C from the predicted values:

$$\ln(y_i) = \beta_0 + \sum_{p=1}^P \beta_p \omega_i^p + \epsilon_i, \quad (4)$$

$$\hat{y}_i^C = \exp \left[\beta_0 + \sum_{p=1}^P \hat{\beta}_p \omega_i^p \right]. \quad (5)$$

The literature refers to this estimation procedure as the parametric approach (Bourguignon et al., 2007; Ferreira and Gignoux, 2011).⁶

According to another procedure, the researcher partitions the sample into mutually exclusive types based on the realizations of all circumstances under consideration. An estimate for y^C is then constructed from average incomes within types:

$$\hat{y}_i^C = \mu_{m(i)} = \frac{1}{N_m} \sum_{j=1}^{N_m} y_j, \quad \forall j \in t_m, \quad \forall t_m \in \mathcal{T}. \quad (6)$$

The literature refers to this estimation procedure as the non-parametric approach (Checchi and Peragine, 2010).

Both approaches face empirical challenges that are typically resolved by discretionary decisions of the researcher. For example, the parametric approach assumes a log–linear impact of all circumstances and therefore neglects the existence of interdependences between circumstances and other non-linearities. To alleviate this shortcoming, researchers can integrate interaction terms and higher-order polynomials into equation (4). However, such extensions remain at their discretion. Reversely, the non-parametric approach does not restrict the interdependent impact of circumstances. However, if the data are rich enough in information on circumstances, researchers might be forced to reduce the observed circumstance vector to obtain statistically meaningful estimates of the relevant parameters.⁷ The

⁶The logarithmic transformation is not innocuous as the marginal impact of circumstances on incomes can differ from their impact on log-incomes. Therefore, the predicted outcome should be obtained by applying the correction suggested in Blackburn (2007). This correction, however, is rarely implemented in empirical applications.

⁷Assume that the researcher observes ten circumstance variables with three expressions each – a quantity easily observed in many datasets. The non-parametric approach would require the estimation of $3^{10} = 59,049$ group means.

8 Estimating inequality of opportunity from regression trees and forests

necessary process of restricting the circumstance vector again remains at the researcher's discretion.

The above discussion illustrates that common approaches leave researchers to their own devices when selecting the best model for estimating the distribution y^C . In this paper, we provide an automated solution to this problem. Similarly, Li Donni et al. (2015) propose the use of latent class modeling to obtain type partitions that allow for estimates of y^C according to the non-parametric procedure outlined in equation (6). In their approach, observable circumstances are considered indicators of membership in an unobservable latent type. For each possible number of latent types, individuals are assigned to types so as to minimize the within-type correlation of observable circumstances. Then the optimal number of types, M^* , is selected by minimizing an appropriate model selection criterion such as Schwarz's Bayesian Information Criterion (BIC). The latent class approach therefore partly solves the issue of arbitrary model selection. However, it has important drawbacks. First, it cannot solve the problem of model selection once the potential number of types exceeds the available degrees of freedom. In such cases, the latent class approach replicates the limitations of parametric and non-parametric approaches: researchers must pre-select circumstances and their subpartition. Second, latent classes are obtained by minimizing the within-type correlation of circumstances while ignoring the correlation of circumstance variables with the outcome variable. As a consequence, they are likely to underfit the data, leading to downward-biased estimates of inequality of opportunity (Lanza et al., 2013).

In the following section, we discuss how regression trees and forests address the outlined shortcomings of existing estimation approaches.

3. Estimating inequality of opportunity from regression trees and forests

Regression trees and forests belong to the class of supervised learning methods that were developed to make out-of-sample predictions of a dependent variable based on a number of observable predictors. As we outline in the following, they can be straightforwardly applied to inequality of opportunity estimations, and they solve the issue of model selection.

First, we introduce conditional inference regression trees. By providing predictions based on identifiable groups, they closely connect to Roemer's theoretical formulation of inequality of opportunity.⁸ Second, we introduce

⁸Furthermore, their simple graphical illustration can be an instructive tool for comparisons of opportunity structures in different societies.

conditional inference forests, which are, loosely speaking, a collection of many conditional inference trees. While forests do not have the intuitive appeal of regression trees, they perform better in terms of out-of-sample prediction accuracy, and hence provide better estimates of the counterfactual distribution y^C .

3.1. Conditional inference trees

Trees obtain predictions for outcome y as a function of input variables $x = (x^1, \dots, x^k)$. They use the sample $\mathcal{S} = \{(y_i, x_i)\}_{i=1}^S$ to divide the population into non-overlapping groups, $\mathcal{G} = \{g_1, \dots, g_m, \dots, g_M\}$, where each group g_m is homogeneous in the expression of some input variables. These groups are called *terminal nodes* or *leafs*. The conditional expectation for observation i is estimated from the mean outcome $\hat{\mu}_m$ of the group g_m to which i is assigned. Hence, in addition to the observed outcome vector $y = (y_1, \dots, y_i, \dots, y_N)$ one obtains a vector of predicted values $\hat{y} = (\hat{f}(x_1), \dots, \hat{f}(x_i), \dots, \hat{f}(x_N))$, where

$$\hat{f}(x_i) = \hat{\mu}_{m(i)} = \frac{1}{N_m} \sum_{j \in g_m} y_j. \tag{7}$$

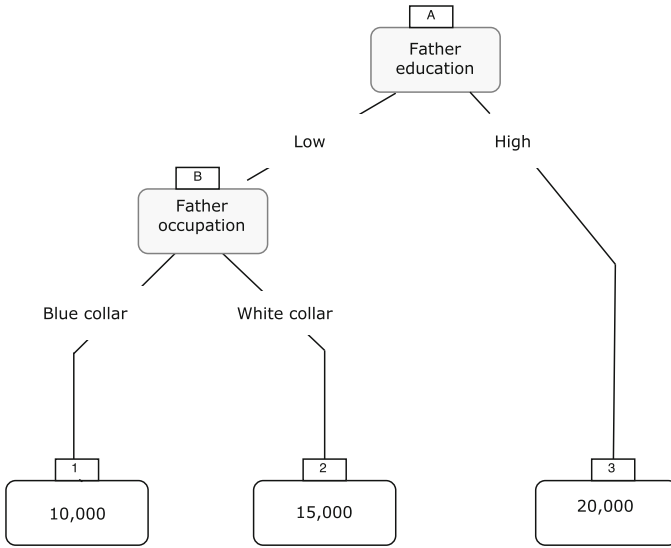
The mapping from regression trees to equality of opportunity estimation is straightforward. If the input variables $x = (x^1, \dots, x^k)$ are circumstances only, each resulting group $g_m \in \mathcal{G}$ can be interpreted as a circumstance type $t_m \in \mathcal{T}$. Furthermore, \hat{y} is analogous to an estimate of the counterfactual distribution y^C that underpins the construction of *ex ante* utilitarian measures of inequality of opportunity.

3.2. Tree construction

Regression trees partition the sample into M types by recursive binary splitting, which starts by dividing the full sample into two distinct groups according to the value they take in one input variable $\omega^p \in \Omega$. If ω^p is a continuous or ordered variable, then $i \in t_l$ if $\omega_i^p < \tilde{\omega}^p$ and $i \in t_m$ if $\omega_i^p \geq \tilde{\omega}^p$, where $\tilde{\omega}^p$ is a splitting value chosen by the algorithm. If ω^p is a categorical variable, then the categories can be split into any two arbitrary groups. The process is continued such that one of the two groups is divided into further subgroups (potentially based on another $\omega^q \in \Omega$), and so on. Graphically, this division into groups can be presented like an upside-down tree (Figure 1).

The exact manner in which the split is conducted depends on the type of regression tree that is used. In this paper, we follow the conditional inference methodology proposed by Hothorn et al. (2006). Conditional inference trees are grown by a series of permutation tests according to the following four-step procedure.

Figure 1. Exemplary tree representation



Notes: Artificial example of a regression tree. Gray boxes indicate splitting points; white boxes indicate terminal nodes. The values inside terminal nodes show estimates for the conditional expectation y^C .

0. Choose a significance level α^* .
1. Test the null hypothesis of density function independence: $H_0^{\omega^p} : D(y|\omega^p) = D(y)$, for all $\omega^p \in \Omega$, and obtain a p -value associated with each test, p^{ω^p} .
 - \Rightarrow Adjust the p -values for multiple hypothesis testing, such that $p_{adj.}^{\omega^p} = 1 - (1 - p^{\omega^p})^P$ (Bonferroni correction).
2. Select the variable ω^* with the lowest p -value, i.e., $\omega^* = \arg \min_{\omega^p} \{p_{adj.}^{\omega^p} : \omega^p \in \Omega, p = 1, \dots, P\}$.
 - \Rightarrow If $p_{adj.}^{\omega^*} > \alpha^*$, exit the algorithm.
 - \Rightarrow If $p_{adj.}^{\omega^*} \leq \alpha^*$, continue, and select ω^* as the splitting variable.
3. Test the null hypothesis of density function independence between the subsamples for each possible binary partition splitting point s based on ω^* , and obtain a p -value associated with each test, $p^{\omega_s^*}$.
 - \Rightarrow Split the sample based on ω^* , by choosing the splitting point s that yields the lowest p -value, i.e., $\tilde{\omega}^* = \arg \min_{\omega_s^*} \{p^{\omega_s^*} : \omega_s^* \in \Omega\}$.
4. Repeat steps 1–3 for each of the resulting subsamples.

In words, conditional inference trees start by a series of univariate hypothesis tests. The circumstance that is most related to the outcome is chosen as the potential splitting variable. If the dependence between the outcome and the splitting variable is sufficiently strong, then a split is made. If not, no split is made. Whenever a circumstance can be split in several ways, the sample is split into two subsamples such that the dependence with the outcome variable is maximized. This procedure is repeated in each of the two subsamples until no circumstance in any subsample is sufficiently related to the outcome variable. Note that the depth of the resulting opportunity tree hinges on the level of α^* . The less stringent the α^* -requirement, the more we allow for false positives (i.e., the more splits will be detected as significant and the deeper the tree will be grown). In our empirical application, we fix $\alpha^* = 0.01$, which is in line with the disciplinary convention for hypothesis tests. To illustrate the robustness of this choice, we show comparisons to setting $\alpha^* = 0.05$ and choosing α^* through cross-validation in Figure S.1 in the Supplementary Material.

A particular advantage of trees is that they avoid list-wise deletion of observations by implementing surrogate splits. In case of missing data, the algorithm searches for an alternative splitting point that mimics the sample partition based on $\tilde{\omega}^*$ to the greatest extent. All observations that lack information on $\tilde{\omega}^*$ are then allocated to subbranches based on this surrogate splitting point.

3.3. Conditional inference forests

Regression trees provide a simple and standardized way of dividing the population into types. Therefore, they solve the model selection problem outlined in Section 2. However, trees suffer from three shortcomings. First, the structure of trees – and therefore the estimate of y^C – is fairly sensitive to alternations in data samples. This issue is particularly pronounced if there are various circumstances that are close competitors for defining the first splits (Friedman et al., 2009). Second, trees assume a non-linear DGP that imposes interactions while ruling out the linear influence of circumstances. Third, trees make inefficient use of data because some of the circumstances $\omega^p \in \Omega$ are not used for the construction of the tree. However, circumstances might possess informational content that can increase predictive power even if they are not significantly associated with y at level α^* . This becomes an issue if two or more important circumstances are highly correlated. Once a split is made using either of the two, it is unlikely that the other contains enough information to cause another split. Conditional inference forests address all of these shortcomings (Breiman, 2001; Biau and Scornet, 2016).

3.4. Forest construction

Random forests create many trees and average over all of these when making predictions. Trees are constructed according to the same four-step procedure outlined in Section 3.2. However, two tweaks are made. First, given the sample $\mathcal{S} = \{(y_i, \omega_i)\}_{i=1}^S$, each tree is estimated on a random subsample $\mathcal{S}' \subset \mathcal{S}$. In our application, we randomly select approximately 60 percent of the observations for each tree, and estimate B^* such trees in total. Second, only a random subset of circumstances of cardinality \bar{P}^* is allowed to be used at each splitting point. Together, these two tweaks remedy the shortcomings of single conditional inference trees. First, averaging over B^* predictions cushions variance in the estimates of y^C and smooths the non-linear impact of circumstance characteristics. Second, the use of subsets of all circumstance variables increases the likelihood that all observed circumstances with informational content will be identified as splitting variable ω^* at some point.

Predictions are formed as follows:

$$\hat{f}(\omega; \alpha^*, \bar{P}^*, B^*) = \frac{1}{B^*} \sum_{b=1}^{B^*} \hat{f}^b(\omega; \alpha^*, \bar{P}^*). \quad (8)$$

Equation (8) illustrates that individual predictions are a function of α^* (i.e., the significance level governing the implementation of splits, \bar{P}^*), the number of circumstances to be considered at each splitting point, and the number of subsamples drawn from the data, B^* . In our empirical illustration, we fix $B^* = 200$ and determine α^* and \bar{P}^* by minimizing the “out-of-bag” error (MSE^{OBB}). Details on these choices and empirical procedures are disclosed in Section S.1 of the Supplementary Material.

4. Empirical application

In this section, we illustrate the machine learning approach using harmonized survey data from 31 European countries. We compare the results from trees and forests with results from the prevalent estimation approaches in the extant literature: parametric, non-parametric, and latent class models. Comparisons are made along two dimensions.

First, we evaluate the different estimation approaches by comparing their out-of-sample mean squared error, MSE^{Test} , which is a standard statistic to evaluate the prediction quality of estimation models.⁹ To calculate MSE^{Test} , we follow the machine learning practice of splitting our sample into a training

⁹Minimizing MSE^{Test} is equivalent to trading off upward and downward biases of inequality of opportunity estimates in a given data environment: the more parsimonious the model, the higher

set with $i^{-H} \in \{1, \dots, N^{-H}\}$ and a test set with $i^H \in \{1, \dots, N^H\}$. For each sample, we choose $N^{-H} = (2/3)N$ and $N^H = (1/3)N$.¹⁰ We fit our models on the training set and compare their performance on the test set according to the following procedure.

1. Run the model on the training data (for the specific estimation procedures, see Section 3.1 for trees and forests, and Section 4.2 for our benchmark methods).
2. Store the prediction function $\hat{f}^{-H}(\cdot)$.
3. Calculate the mean squared error in the test set:

$$\text{MSE}^{\text{Test}} = (1/N^H) \sum_{i \in H} [y_i - \hat{f}^{-H}(\omega_i)]^2.$$

Second, we evaluate the different approaches by comparing inequality of opportunity estimates. To this end, we run the models on all data for a country, and apply the resulting prediction functions $\hat{f}(\cdot)$ to obtain \hat{y}^C . Estimates of inequality of opportunity are derived by summarizing \hat{y}^C with the Gini index. Estimates for alternative inequality indices are presented in Section S.6 of the Supplementary Material.

4.1. Data

We base our empirical illustration on the 2011 wave of the European Union Statistics on Income and Living Conditions (EU-SILC), which provides harmonized survey data with respect to income, poverty, and living conditions. It is the official reference source for comparative statistics on income distribution and social inclusion in the EU. In its 2011 wave, EU-SILC covers a cross-section of 31 European countries. For each country, it contains a random sample of all resident private households. Data are collected by national statistical agencies following common variable definitions and data collection procedures. We use the 2011 wave because it contains an ad hoc module about the intergenerational transmission of (dis)advantages. This module allows us to construct finely grained circumstance-type partitions. Observed circumstances Ω and their respective expressions are listed in Table 1. We include all variables of EU-SILC containing information about

the prediction bias (underfitting) and the stronger the downward bias in inequality of opportunity estimates. The more complex the model, the higher the prediction variance (overfitting) and the stronger the upward bias of inequality of opportunity estimates. We provide a thorough illustration of this mapping in Section S.2 in the Supplementary Material.

¹⁰Note that the size of the training set for each country is constant regardless of the estimation method. Hence, any cross-method differences in prediction accuracy are not driven by differences in sample size.

14 Estimating inequality of opportunity from regression trees and forests

Table 1. List of circumstances

-
- (1) Respondent's sex: male; female
 - (2) Respondent's country of birth:
present country of residence; European country; non-European country
 - (3) Presence of parents at home*:
both present
only mother
only father
without parents
lived in a private household without any parent
 - (4) Number of adults (aged 18 or more) in respondent's household*
 - (5) Number of working adults (aged 18 or more) in respondent's household*
 - (6) Number of children (under 18) in respondent's household*
 - (7) Father's/mother's country of birth and citizenship:
born in/citizen of the respondent's present country of residence
born in/citizen of another EU-27 country
born in/citizen of another European country
born in/citizen of a country outside Europe
 - (8) Father's/mother's education (based on ISCED-97)*:
unknown father/mother
illiterate
low (0–2 ISCED-97), medium (3–4 ISCED-97) or high (5–6 ISCED-97)
 - (9) Father's/mother's occupational status*:
unknown or dead father/mother
employed
self-employed
unemployed
retired
house worker
other inactive
 - (10) Father's/mother's main occupation (based on ISCO-08)*:
managers (I-01)
professionals (I-02)
technicians (I-03)
clerical support workers (I-04)
service and sales workers (I-05 and 10)
skilled agricultural, forestry and fishery workers (I-06)
craft and related trades workers (I-07)
plant and machine operators, and assemblers (I-08)
elementary occupations (I-09)
armed forces occupations (I-00)
father/mother did not work, was unknown or was dead
 - (11) Managerial position of father/mother*: supervisory; non-supervisory
 - (12) Tenancy status of the house in which the respondent was living*: owned; not owned
-

Notes: This table lists the circumstance variables available in EU-SILC 2011. Questions marked with * refer to the time when the respondent was 14 years old. Item 7 (11) is missing for Slovenia (Finland). ISCED97 is the International Standard Classification of Education 1997. ISCO-08 is the International Standard Classification of Occupations, published by the International Labour Office.

Source: EU-SILC 2011 cross-sectional (rev. 5, June 2015).

the respondent's characteristics at birth and their living conditions during childhood. Descriptive statistics of circumstance variables are reported in Section S.5 of the Supplementary Material.¹¹

The unit of observation is the individual and the outcome of interest is equivalized disposable household income. We obtain the latter by dividing household disposable income with the square root of household size. Reported incomes refer to the year preceding the survey wave (i.e., 2010 in the case of our empirical application). In line with the literature, we focus on equivalized household income as it provides the closest income analogue to consumption possibilities and general economic well-being. Inequality statistics tend to be heavily influenced by outliers (Cowell and Victoria-Feser, 1996); therefore, we adopt a standard winsorization method according to which we set all non-positive incomes to 1 and scale back all incomes exceeding the 99.5th percentile of the country-specific income distribution to this lower threshold. Our analysis is focused on the working-age population. Therefore, we restrict the sample to respondents aged between 30 and 60. To assure the representativeness of our inequality of opportunity estimates, we use individual cross-sectional weights when calculating $I(\hat{y}^C)$.

Table 2 shows considerable heterogeneity in income distributions across Europe. While the average households in Norway and Switzerland obtained incomes above 40,000 euros in 2010, the average household income in Romania, Bulgaria, and Lithuania did not exceed the 5,000 euros mark. Lowest inequality prevails in Norway, Iceland, and Denmark, all of which have Gini coefficients below 0.230. At the other end of the spectrum, we find Latvia and Lithuania with Gini coefficients above 0.340.

4.2. Benchmark methods

We compare trees and forests to three benchmark estimation methods from the extant literature.

First, we draw on the parametric approach as proposed by Bourguignon et al. (2007) and Ferreira and Gignoux (2011). In line with equation (4), estimates are obtained by a Mincerian regression of log income on the following circumstances: educational attainment of mother and father (five categories each), father's occupation (11 categories), area of birth (three categories), and tenancy status of the household at age 14 (two categories). The prediction model includes 22 parameters.

¹¹In contrast to some existing work, we do not consider age as a circumstance (see Checchi et al., 2016, among others). This choice is motivated by the fact that cross-sectional income disparities across age groups even out across the life cycle of individuals. In Section S.7 of the Supplementary Material, we provide robustness analyses based on income distributions that are residualized from variation across age groups. Our conclusions remain unaffected.

Table 2. Summary statistics

Country	<i>N</i>	Equivalent disposable household income in euros		
		μ	σ	Gini
Austria	6,220	25,538	13,408	0.267
Belgium	6,011	23,314	10,769	0.247
Bulgaria	7,146	3,698	2,457	0.331
Croatia	6,945	6,631	3,764	0.306
Cyprus	4,589	21,074	11,554	0.278
Czech Republic	8,711	9,036	4,610	0.253
Denmark	5,795	32,471	14,422	0.227
Estonia	5,338	6,924	4,364	0.330
France	11,078	24,320	14,695	0.287
Germany	12,683	22,862	12,468	0.284
Greece	6,184	13,184	8,887	0.334
Hungary	13,330	5,305	2,830	0.275
Iceland	3,682	21,562	9,290	0.221
Ireland	4,318	24,882	14,078	0.295
Italy	21,070	18,774	11,348	0.314
Latvia	6,423	5,339	3,751	0.362
Lithuania	5,403	4,774	3,116	0.344
Luxembourg	6,765	37,948	19,412	0.270
Malta	4,701	13,058	6,758	0.272
Netherlands	11,411	24,322	11,452	0.243
Norway	5,026	42,265	16,679	0.206
Poland	15,545	6,087	3,837	0.316
Portugal	5,899	10,796	7,354	0.333
Romania	7,820	2,527	1,612	0.336
Slovakia	6,779	7,309	3,509	0.256
Slovenia	13,183	13,373	5,896	0.234
Spain	15,481	17,088	10,684	0.328
Sweden	6,599	25,098	11,157	0.237
Switzerland	7,583	42,844	23,877	0.278
United Kingdom	7,391	22,768	15,164	0.319

Notes: This table provides summary statistics by country. *N* indicates the total number of observations. The last three columns summarize the distribution of equivalized disposable household income: mean (μ), standard deviation (σ), and Gini coefficient.

Source: EU-SILC 2011 cross-sectional (rev. 5, June 2015).

Second, we draw on the non-parametric approach as proposed by Checchi and Peragine (2010). In line with equation (6), non-parametric estimates are obtained by calculating average outcomes in non-overlapping circumstance types. Types are homogeneous with respect to educational attainment of the highest educated parent (five categories), fathers' occupation (four categories),

and migration status (two categories).¹² The prediction model includes 40 parameters.

Third, we draw on the latent class approach as proposed by Li Donni et al. (2015). We use the union of circumstances used in the parametric and non-parametric approaches from which the algorithm infers the appropriate number of unobserved types in the data by minimizing the BIC.

Do these specification choices serve for a fair assessment of these benchmark methods? As outlined in Section 2, model specification in the (non-)parametric approach is a discretionary choice of the researcher; therefore, there are many different specifications that could be used for the benchmarking. To make the comparison non-arbitrary, we anchor our comparison on model specifications of existing studies. The specification of the parametric approach is inspired by Palomino et al. (2019). We divert from their specification by excluding gender (due to our focus on disposable household income) and retrospective information on the financial situation during childhood (due to potential recall bias) from the list of circumstances. In comparison, our prediction model (22 parameters) is more parsimonious than the model in Palomino et al. (2019, 24 parameters). The specification of the non-parametric approach is inspired by Checchi et al. (2016). We divert from their specification by excluding gender (due to our focus on disposable household income) and age (due to its interpretation as a proxy for life-cycle effects) from the list of circumstances. In comparison, our prediction model (40 parameters) is more parsimonious than the model in Checchi et al. (2016, 96 parameters). As outlined in Section 2, model specification in latent class analysis is data-driven. Therefore, we do not need to specify the model itself but commit to a model selection criterion. We anchor our comparison on the study of Li Donni et al. (2015) who use the BIC to select the number of latent classes to be estimated.

4.3. Simulation

We begin our analysis with a simulation exercise. The simulation allows us to assess the properties of different estimation approaches while maintaining control over the true DGP. As a consequence, we can: (i) assess the prediction accuracy by decomposing MSE^{Test} into its variance and bias components; and (ii) assess the resulting bias in inequality of opportunity estimates.

In general, simulation results are sensitive to assumptions about the true DGP and sample sizes. To make the simulation relevant to the context of

¹²To minimize the frequency of sparsely populated types, we divert from the occupational list given in Table 1 by re-coding occupations into the following categories: high-skilled non-manual (I-01–I-03), low-skilled non-manual (I-04–I-05 and I-10), skilled manual and elementary occupation (I-06–I-09), and unemployed/unknown/dead.

Table 3. Summary of data-generating processes

	Parametric	Non-parametric	Mixed
Outcome	$\ln(y)$	y	$\ln(y)$
Parameters	22	40	18
Circumstances	Education father	Education parents	Education parents
	Education mother	Occupation father	Occupation father
	Occupation father	Migrant background	Migrant background
	Birth area		Tenancy status
	Tenancy status		
Non-linearity	None	Full interaction	All circumstances with migrant background (two levels)
ϵ	$\mathcal{N}(0, 2000)$	$\mathcal{N}(0, 2000)$	$\mathcal{N}(0, 2000)$

our empirical analysis, we choose DGPs that are anchored in the benchmark methods presented in Section 4.2 and choose sample sizes to broadly cover the range of country samples in EU-SILC. We note that our additional simulation choices are conservative. First, we construct a simulation sample without missing data points. As a consequence – and in contrast to actual empirical applications – parametric and non-parametric approaches do not suffer from data reductions through list-wise deletion. Second, we restrict circumstances used by trees and forests to the union of circumstances used in the (non-)parametric approach. As a consequence – and in contrast to actual empirical applications – we deprive data-driven approaches from the advantage of using all available circumstance information in the data.

We impose three DGPs that are summarized in Table 3. The parametric DGP and the non-parametric DGP correspond to the estimation models outlined in Section 4.2. They present a challenging test for data-driven estimation methods because the latter have to compete against fixed specifications (parametric, non-parametric) that correspond to the ground truth. In addition, we specify a mixed DGP that integrates features of both the parametric and the non-parametric DGP. This is a more realistic scenario as it is plausible to assume that researchers devise fixed specifications without prior knowledge of the true DGP. We estimate all three models on the full sample of EU-SILC while list-wise deleting observations with missing information ($N = 197,565$). In turn, we retain the predictions from these estimations and add a disturbance term with $\mathcal{N}(0, 2000)$.¹³ Thus, we obtain three variables that define the distribution of income for the purpose of this simulation.

¹³We choose a variance term small enough such that $y \in \mathbb{R}_{++}$.

Next we specify five sample sizes that broadly cover the range of effective country sample sizes observed in EU-SILC (see Table S.1 in the Supplementary Material): $N \in \{1,000; 2,000; 4,000; 8,000; 16,000\}$. For each sample size, we draw one test set of size $N^H = (1/3)N$ and 50 training sets of size $N^{-H} = (2/3)N$. Thus, for each observation in the test sets, we obtain 50 predictions per combination of DGP and estimation approach. Based on these predictions, we calculate two statistics: the expected MSE^{Test} to assess out-of-sample prediction accuracy (James et al., 2013), and the expected absolute difference between inequality of opportunity estimates and the true level of inequality of opportunity.

Figure 2 displays the results of our simulation. The lower part of each panel describes expected MSE^{Test} per combination of DGP, estimation approach, and sample size. As we know the true DGP, we can decompose MSE^{Test} into variance and expected bias.¹⁴ The upper part of each panel describes the corresponding absolute bias in inequality of opportunity estimates on an inverse scale. The absolute bias is calculated as the expected absolute difference between inequality of opportunity estimates and the true level of inequality of opportunity, as a percentage share of the latter.

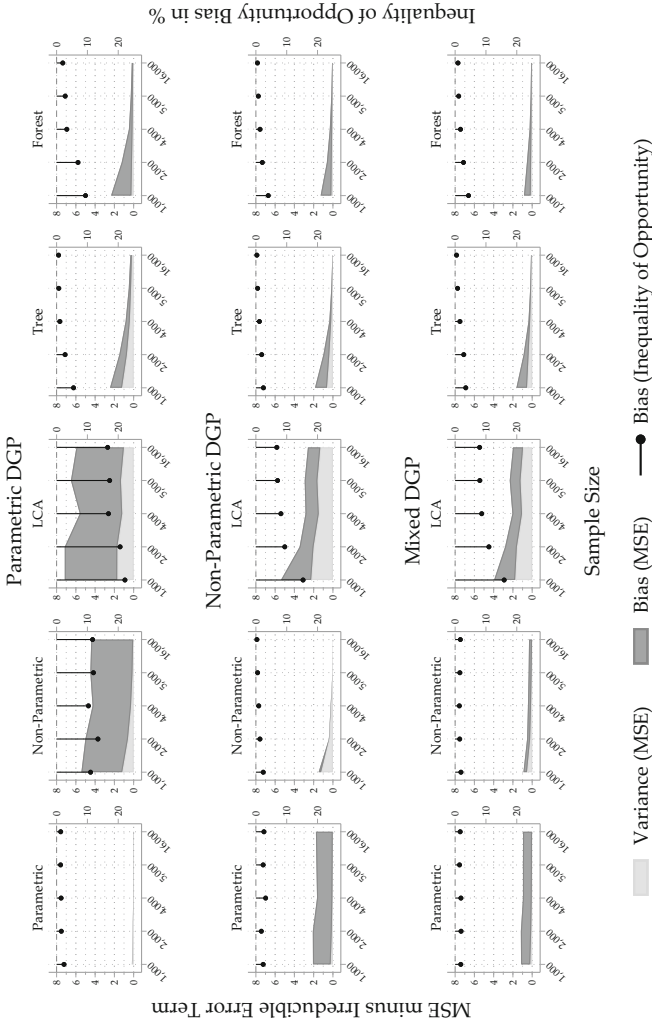
The simulation results are in line with statistical theory. First, if fixed estimation approaches (parametric, non-parametric) invoke the true DGP, expected bias is zero and MSE^{Test} is driven by its variance component only. Second, with increasing N , the bias component of MSE^{Test} remains constant for fixed specifications (parametric, non-parametric) and decreases for data-driven approaches (LCA, trees, forests). Third, with increasing N , the variance component of MSE^{Test} decreases for all combinations of DGPs and estimation approaches. Fourth, forests tend to have lower variance than trees – in our simulation, this is true in 93 percent of all cases.

Furthermore, Figure 2 illustrates that the size of MSE^{Test} , and therefore bias in inequality of opportunity estimates, varies with sample size for all estimation methods. However, the sources of this bias vary across estimation methods. For example, forests incur downward bias in small samples as their algorithm prevents the detection of relevant splits. To the contrary, fixed specifications incur upward bias in small samples as the underlying parameters are noisily estimated (although unbiased in expectation). The crucial question is whether this sensitivity is larger for machine learning than for traditional econometric methods, which is a case-specific and empirical question.

Under the reasonable assumption that researchers do not know the true DGP, forests clearly dominate all other estimation approaches in terms of

¹⁴See also Section S.2 in the Supplementary Material for an illustration of the variance–bias decomposition. The irreducible error term is uninformative for differences in MSE^{Test} because $\text{Var}(\epsilon) = 2,000^2$ is constant across specifications. Therefore, we only present evidence on the variance and the bias component of MSE^{Test} .

Figure 2. Simulation results



Notes: This figure shows expected MSE^{Test} and expected bias in inequality of opportunity estimates from the simulation exercise. Each row corresponds to one DGP (see Table 3). Each column corresponds to one estimation method (see Sections 3.1, 3.3, and 4.2). We multiply MSE^{Test} by 1×10^{-6} and deduct the irreducible error term. Inequality of opportunity is measured by the Gini coefficient of the counterfactual distribution \hat{y}^C . We measure expected bias in inequality of opportunity by the average absolute difference between inequality of opportunity estimates and the true level of inequality of opportunity as specified by the DGP.

Source: EU-SILC 2011 cross-sectional (rev. 5, June 2015).

expected MSE^{Test} . This result holds both in comparison with fixed estimation approaches (parametric, non-parametric) and in comparison with LCA as an alternative data-driven estimation approach. The results for trees are also persuasive; however, they have a weaker performance than forests when samples are small. Even in the unlikely case that researchers were to specify (non-)parametric models correctly, trees and forests quickly converge to the test error of the fixed model that invokes the true DGP. The simulation results further highlight the close correspondence between MSE^{Test} and expected bias in inequality of opportunity estimates: the higher MSE^{Test} , the more strongly inequality of opportunity estimates diverge from the ground truth.

In summary, the simulation results support the use of regression trees and forests. They flexibly approximate the true DGP. Thereby, they outperform fixed estimation approaches (parametric, non-parametric) and alternative data-driven estimation approaches (LCA) in terms of the expected MSE^{Test} , which itself is tightly linked to expected bias in inequality of opportunity estimates.

4.4. Cross-country comparison

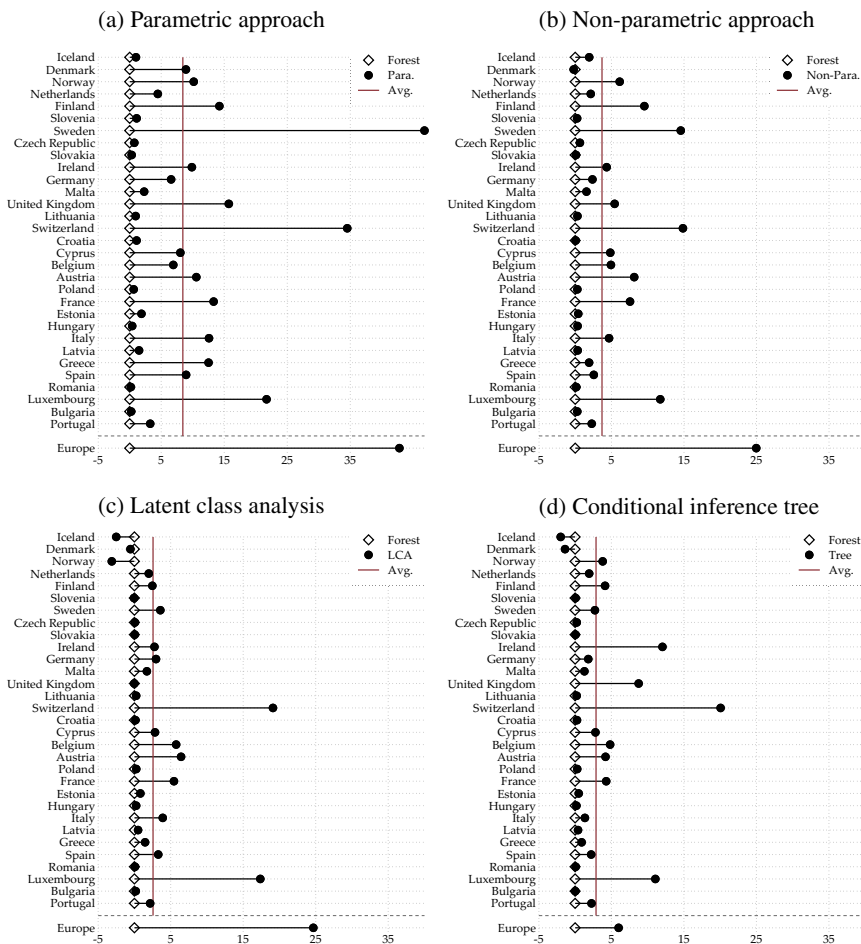
We now turn to a cross-country comparison based on actual data. First, we calculate MSE^{Test} to assess the prediction accuracy of different estimation approaches. Second, we calculate inequality of opportunity estimates. In contrast to the simulation exercise, we do not know the true DGP and we cannot assess bias in inequality of opportunity estimates by comparison with the ground truth. Therefore, we assess bias in inequality of opportunity estimates by comparing estimation approaches against the method with the highest prediction accuracy (i.e., the method yielding the lowest MSE^{Test}).

4.4.1. Prediction accuracy. Figure 3 compares MSE^{Test} across countries and estimation approaches. For each method, MSE^{Test} is presented in differences relative to random forests. By differencing across methods, we provide a close analogue to the simulation exercise in Section 4.3: we omit the irreducible error term from the comparison, and relative MSE^{Test} is driven by variance and bias components only. For better visual clarity, we again scale MSE^{Test} by 1×10^{-6} . Relative $MSE^{Test} > 0$ indicates poorer prediction accuracy in comparison with random forests.

Random forests outperform all other methods in terms of prediction accuracy. On average, the parametric approach yields test errors that exceed random forests by 8.4 (7.8 percent); see Figure 3(a). Somewhat smaller average shortfalls in prediction accuracy are observed for non-parametric (Figure 3(b)) and latent class models (Figure 3(c)). Averages across countries,

22 Estimating inequality of opportunity from regression trees and forests

Figure 3. Comparison of MSE^{Test} by method



Note: This figure shows differences of MSE^{Test} from different estimation approaches relative to random forests. For all methods, we multiply MSE^{Test} by 1×10^{-6} . Values $>$ ($<$)0 indicate worse (better) out-of-sample prediction accuracy than random forests. Vertical lines indicate unweighted cross-country averages. Point estimates and associated standard errors are listed in Table S.1 in the Supplementary Material.

Source: EU-SILC 2011 cross-sectional (rev. 5, June 2015).

however, mask considerable heterogeneity. For example, the relative test error of parametric estimates for Eastern European countries, such as Slovenia or Czech Republic, are close to zero. On the contrary, relative test errors of parametric estimates for Sweden, Luxembourg, and Switzerland diverge significantly from the forest benchmark.

Conditional inference trees are closest to the test error rate of forests: $MSE^{\text{Test}} = 2.8$ (2.2 percent). Yet, they also fall short of the performance of forests due to higher variance, imposing non-linearity, and omitting less relevant circumstances (see Section 3.3).

We conclude that among all considered methods, conditional inference forests deliver the highest out-of-sample prediction accuracy.¹⁵ Hence, relative to random forests, other methods underutilize or overutilize the information contained in Ω , which in expectation will lead to bias in inequality of opportunity estimates.

4.4.2. Inequality of opportunity estimates. Figure 4 displays inequality of opportunity estimates across countries and estimation approaches. In each panel, we plot inequality of opportunity estimates for a particular method, as well as the associated differences to estimates from forests. We emphasize that results from forests cannot be interpreted as the truth. However, because forests yield the lowest test error among all considered methods, they provide the best “approximation of the true DGP in a given data environment”. Therefore, they are a useful benchmark to assess bias of other estimation methods.¹⁶

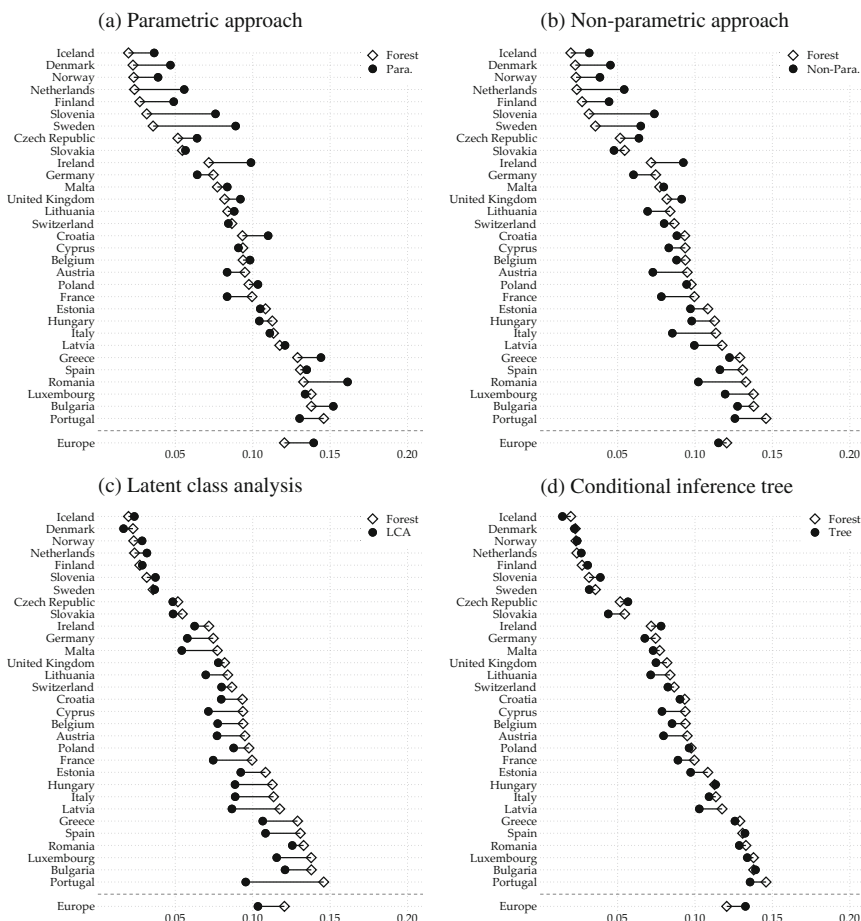
Figure 4(a) shows estimates from the parametric approach. In our country sample, the chosen model specification for the parametric approach tends to overstate inequality of opportunity relative to forests, which is the method providing the lowest expected bias in comparison with the true DGP. For 21 out of 31 countries, the inequality of opportunity estimates are higher than the results from forests. Most pronounced overstatements are observed in countries that are typically considered high-opportunity societies. For example, forests classify Sweden and the Netherlands as societies offering high equality of opportunity. On the contrary, the parametric estimate would rank them at similar levels to Germany and France.

Figure 4(b) shows estimates from the non-parametric approach. The overall pattern is more heterogeneous than for the parametric approach.

¹⁵In Table S.3 of the Supplementary Material, we show that the overwhelming majority of detected differences are statistically significant at conventional levels.

¹⁶It is important to keep this relative interpretation of “bias” in mind. We compare method-specific estimates to the best estimate of inequality of opportunity in a given data environment. In light of all methods lacking information on unobserved circumstances, methods that are upward biased in this comparison might potentially be closer to the ground truth than our reference estimate. Such conclusions, however, are purely speculative and can neither be verified nor falsified without knowledge of the ground truth (see also Section S.2 of the Supplementary Material for a thorough explanation). Therefore, another interpretation of forests is that they provide the reliable maximum lower bound estimate of inequality of opportunity in a given data environment.

Figure 4. Comparison of inequality of opportunity estimates by method



Notes: This figure shows inequality of opportunity estimates from different estimation methods relative to forests. Inequality of opportunity is measured by the Gini coefficient of the counterfactual distribution y^C . Point estimates and associated standard errors are listed in Table S.2 of the Supplementary Material.

Source: EU-SILC 2011 cross-sectional (rev. 5, June 2015).

While overstatements prevail in countries that are typically considered as high-opportunity societies, there are 20 out of 31 countries for which the non-parametric estimate falls short of the forest estimate. These countries are clustered in the lower part of the equal-opportunity ranking. For example, forests classify Italy at a worse position than most countries in Europe. On the contrary, the non-parametric estimate would elevate Italy towards the mid-field, close to Sweden.

We highlight that any resemblance between forests and (non-)parametric estimation approaches should be interpreted as a luck of the draw rather than a property inherent to the estimation approach. Under alternative plausible model specifications, estimates from both approaches might diverge much more strongly than under the specifications adopted in this work. This property of fixed model specifications is apparent from the simulation results in Section 4.3.

Figure 4(c) shows estimates from the latent class approach. In our country sample, LCA tends to understate inequality of opportunity relative to forests. For 25 out of 31 countries, the LCA estimate falls short of its forest-based analogue. LCA chooses rather coarse type partitions. Therefore, understatements are clustered in the lower tail of the equal-opportunity ranking (i.e., in societies in which many circumstances co-produce the outcome of interest). On the contrary, in high-opportunity societies, the parsimonious models chosen by LCA tend to replicate the results from forests reasonably well.

Figure 4(d) shows that trees and forests tend to produce similar results. The correlation between point estimates is high (0.98). In contrast to all other approaches, there is no general tendency to overestimate or underestimate inequality of opportunity relative to forests.

Finally, detected differences between the benchmark estimation approaches and forests persist when estimating equality of opportunity in a pooled European sample.¹⁷ For example, the parametric approach overestimates inequality of opportunity relative to forests, whereas LCA yields lower estimates than forests.

4.4.3. Robustness to differences in sample size. Effective sample sizes differ by estimation method and country (see Table S.1 in the Supplementary Material). First, samples for the benchmark methods (parametric, non-parametric, LCA) are reduced as they rely on list-wise deletion in case of missing circumstance information. These reductions can be sizable and exceed 50 percent in six countries of our sample (Denmark, Iceland, Netherlands, Norway, Slovenia, and Sweden). Second, even when accounting for missing information, the largest country sample in EU-SILC (Italy, $N = 21,070$) is almost seven times as large as the smallest country sample (Iceland, $N = 3,682$). Therefore, we perform two robustness analyses.

¹⁷Note that we do not include country of residence as a circumstance. We acknowledge that country of residence is congruent with country of birth for most individuals. Therefore, it could be used as a proxy circumstance (Milanovic, 2015). However, our foremost concern is a methodological comparison of estimation approaches in different data environments. Therefore, we prefer to keep the set of circumstances comparable to our within-country estimates.

First, we recompute inequality of opportunity after completing missing data through multiple imputation (Schafer, 1999).¹⁸ As a consequence, we can compare inequality of opportunity estimates across methods on the same effective sample size per country. Figure S.1 in the Supplementary Material shows a decrease in inequality of opportunity estimates relative to forests for all benchmark methods (parametric, non-parametric, LCA). This result is in line with the intuition that upward biases decrease as sample sizes grow relative to the number of model parameters. In contrast, the patterns for trees and forests remain unaffected as they handle missing values by default through surrogate splits.

Second, we recompute inequality of opportunity while reducing sample sizes across countries to the smallest common denominator. As a consequence, we can compare inequality of opportunity estimates across countries on the same effective sample size. Figure S.2 in the Supplementary Material shows that point estimates and country rankings differ for all benchmark methods (parametric, non-parametric, LCA). Furthermore, trees also show some variability as sample sizes decrease. In contrast, point estimates and country rankings of forests are unaffected by harmonization in sample sizes across countries. On the one hand, these results highlight that the high variance of trees can lead to suboptimal results in some applications and that researchers should give preference to forest estimates where possible. On the other hand, these results bolster confidence that opportunity rankings of forests are not an artifact of cross-country variation in sample sizes.¹⁹

4.4.4. Comparison with existing literature. We have shown that benchmark methods from the existing literature yield markedly different estimates of inequality of opportunity relative to the method for which we expect the lowest bias. These differences are manifested in both point estimates and country rankings. Therefore, these methods can be misleading in two related dimensions. First, they might mis-classify European societies regarding their need for opportunity equalizing policy interventions. Second, researchers and policymakers in search of best practices to devise opportunity equalizing policy interventions might turn to the wrong country examples. In

¹⁸List-wise deletion yields unbiased parameter estimates if data are missing completely at random (MCAR). Multiple imputation weakens this assumption by assuming that data are missing at random (MAR; i.e., missing data are random conditional on observed variables).

¹⁹We perform a similar exercise on the pooled sample: we re-estimate our results for the pooled sample on increasingly smaller fractions of the total sample. In Figure S.3 in the Supplementary Material, we again show that benchmark methods (parametric, non-parametric, LCA) and trees are sensitive to changes in sample size when fractions become small. In contrast, forests again emerge as the method that is most robust in small samples.

the following, we assess the extent to which such concerns are reflected in the extant literature on inequality of opportunity in Europe.

We proceed in two steps. First, we assess whether the existing literature on inequality of opportunity in Europe is consistent (i.e., whether it yields similar opportunity rankings across European societies). If the literature were consistent, then researcher discretion in model selection would be irrelevant for conclusions about inequality of opportunity in Europe. Second, we assess whether the existing literature on inequality of opportunity in Europe conforms with evidence on the intergenerational income elasticity (IGE). The IGE is a commonly used proxy statistic for equality of opportunity that is based on data links across generations. The IGE provides a suitable benchmark as it can be interpreted as an *ex ante* utilitarian measure of inequality of opportunity (see footnote 5) and it is often based on richer (administrative) panel data. If there was conformity, then current estimation approaches would yield opportunity rankings that are strongly in line with common priors about mobility in European societies. We answer both questions by calculating correlations in opportunity rankings across: (i) existing studies on inequality of opportunity,²⁰ (ii) existing consensus estimates of the IGE,²¹ and (iii) inequality of opportunity estimates from our preferred methods (i.e., regression trees and forests).

Panel A of Table 4 suggests that the existing literature on inequality of opportunity in Europe is not consistent. Rank correlations as low as 0.09 indicate strong heterogeneity in country rankings. On the one hand, all studies that inform this comparison have a very high degree of harmonization in relevant dimensions: estimates were derived from the same underlying data source (EU-SILC), refer to a similar age group (ages 25–60), and summarize counterfactual distributions \hat{y}^C by the same inequality metric (mean log deviation). On the other hand, all studies specify different prediction functions to estimate inequality of opportunity. This suggests that discretionary choices with respect to model specifications might be a major force behind inconclusive evidence in the inequality of opportunity literature. We cannot fully rule out the possibility that differences in income concept definitions – that is, individual income (Checchi et al., 2016) versus household income (Palomino

²⁰We focus on published studies estimating *ex ante* measures of inequality of opportunity on the 2011 wave of EU-SILC. Further studies that do not meet both criteria include Andreoli and Fusco (2019) and Carranza (2020). Furthermore, we do not include Brzezinski (2020) as he derives estimates based on the methods proposed in this paper.

²¹We focus on IGE estimates based on actual data linkages across generations and we exclude IGE estimates based on two-sample instrumental variable estimators to mitigate distortions through measurement error. Estimates are extracted from Stuhler (2018) and Carmichael et al. (2020). Jointly both studies contain the following subset of our country sample: Denmark, Finland, France, Germany, Italy, Netherlands, Norway, Sweden, Spain, and the United Kingdom.

Table 4. Rank correlations of existing studies

	Existing studies			This paper	
	Checchi et al. (2016)	Palomino et al. (2016)	Suárez Álvarez and López Menéndez (2021)	Tree	Forest
Panel A. Equality of opportunity (23 countries)					
Tree	–	–	–	1.000	–
Forest	–	–	–	0.984	1.000
Checchi et al. (2016)	1.000	–	–	0.363	0.345
Palomino et al. (2019)	0.281	1.000	–	0.882	0.859
Suárez Álvarez and López Menéndez (2021)	0.090	0.855	1.000	0.756	0.757
Panel B. Intergenerational elasticity (10 countries)					
Stuhler (2018) and Carmichael et al. (2020)	0.535	0.657	0.444	0.900	0.887

Notes: This table shows country rank correlations in inequality of opportunity estimates across existing studies. Panel A is based on the intersection of countries included in this paper, Checchi et al. (2016), Palomino et al. (2019), and Suárez Álvarez and López Menéndez (2021) (23 countries). All ranks are calculated from the mean log deviation of the counterfactual distribution \hat{y}^C . Panel B is based on the intersection of countries included in this paper, Palomino et al. (2019), Checchi et al. (2016), and Suárez Álvarez and López Menéndez (2021), and the union of Stuhler (2018) and Carmichael et al. (2020) (10 countries). Ranks in Stuhler (2018) and Carmichael et al. (2020) are calculated from consensus estimates of the intergenerational earnings elasticity (IGE). All rank correlations are based on Spearman's ρ .

Source: EU-SILC 2011 cross-sectional (rev.5, June 2015).

et al., 2019; Suárez Álvarez and López Menéndez, 2021) – might also contribute to observed lack of consistency. However, as we detail in the next paragraph, regardless of their income definition, all of these studies are only moderately correlated with IGE rankings that are calculated with respect to individual incomes. This pattern does not support an alternative explanation based on differences in income definitions (see Panel B of Table 4).

In Panel B of Table 4, we test for conformity of opportunity rankings with the IGE literature. Inequality of opportunity rankings of existing studies are only moderately correlated with IGE rankings. In fact, various findings contradict comparative evidence on the IGE (Carmichael et al., 2020; Bratberg et al., 2017). For example, Palomino et al. (2019) and Suárez Álvarez and López Menéndez (2021) find inequality of opportunity in Germany to be on a par with the Nordic countries. Checchi et al. (2016) find the Netherlands to be in the lower part of the opportunity ranking. On the contrary, rankings based on trees and forests strongly increase conformity with IGE estimates and therefore yield results that are more strongly in line with common priors about mobility in Europe.

We conclude that regression trees and forests foster consistency in the inequality of opportunity literature by reducing researcher discretion and increase conformity with evidence from the neighboring IGE literature. Both findings further bolster confidence in the ability of trees and forests to make reliable distinctions among high and low opportunity societies in Europe.

5. Conclusion

In this paper we propose the use of conditional inference trees and forests to estimate inequality of opportunity. Both estimation approaches minimize arbitrary model selection by the researcher while trading off downward and upward biases in inequality of opportunity estimates.

Conditional inference forests outperform all methods considered in this paper in terms of their out-of-sample prediction accuracy. This observation is valid both for simulated DGPs and representative survey data from 31 European countries. Hence, within a given data environment, they provide estimates of inequality of opportunity that have the lowest expected bias. Conditional inference trees closely mirror forests in terms of their out-of-sample prediction accuracy and their inequality of opportunity estimates. Hence, they provide a fair first-order approximation to the least-biased inequality of opportunity estimates. Nevertheless, researchers should be conscious that trees might yield suboptimal results in applications with smaller samples.

We note that the improvements of our preferred methods are conditional on a given data environment. As a consequence, they do not address two major challenges of the existing literature: bias due to unobserved circumstance information, and bias due to small sample sizes. These challenges exist independently of the chosen estimation technique and can only be overcome through the availability of improved data sources in the future.

Next to their advantages, we acknowledge two potential drawbacks of our preferred methods for empirical research. First, (non-)parametric estimation approaches can be estimated by ordinary least squares (OLS) – one of the workhorse estimation methods in economics and other social sciences. To the contrary, machine learning tools might require some upfront investment of applied researchers to familiarize themselves with these methods. However, as evidenced by the large volume of recent review articles, machine learning methods are increasingly integrated into the statistical toolkit of economists (Varian, 2014; Mullainathan and Spiess, 2017; Athey, 2018). Therefore, we expect this drawback to vanish over time. Second, trees and forests are computationally more costly than predictions via OLS regressions. However, in our empirical application, trees approach

the computation times of the (non-)parametric approach.²² Therefore, time-constrained researchers who are willing to settle for a fair first-order approximation of the least-biased method can consider using trees instead of forests.

The development of machine learning algorithms and their integration into the analytical toolkit of economists is a dynamic process. Finding the best machine learning algorithm for inequality of opportunity estimations is a methodological horse race that eventually will lead to some method outperforming the ones employed in this work. Therefore, the main contribution of this work should be understood as paving the way for new methods that are able to handle the intricacies of model selection for inequality of opportunity estimations. A particularly interesting extension may be the application of local linear forests that outperform more traditional forest algorithms in their ability to capture the linear impact of predictor variables (Friedberg et al., 2020).

Finally, we restricted ourselves to *ex ante* utilitarian measures of inequality of opportunity. The exploration of these algorithms for other measurement approaches in the inequality of opportunity literature provides another interesting avenue for future research (Lefranc et al., 2009; Pistolesi, 2009; Kanbur and Snell, 2019; Brunori and Neidhöfer, 2021).

Supporting information

Additional supporting information can be found online in the supporting information section at the end of the article.

Supplementary material Replication files

References

- Alesina, A., Stantcheva, S., and Teso, E. (2018), Intergenerational mobility and preferences for redistribution, *American Economic Review* 108 (2), 521–554.
- Andreoli, F. and Fusco, A. (2019), Robust cross-country analysis of inequality of opportunity, *Economics Letters* 182, 86–89.
- Arneson, R. J. (2018), Four conceptions of equal opportunity, *Economic Journal* 128, F152–F173.

²²The simulation of Section 4.3 has the following computation times: 0.5 min (parametric), 0.5 min (non-parametric), 121.4 min (LCA), 2.1 min (trees), 2,816.2 min (forests). These computation times are based on a machine with an AMD Ryzen 7 4700U Processor (8 cores) and 16 GB RAM working memory.

- Athey, S. (2018), The impact of machine learning on economics, in A. K. Agrawal, J. Gans and A. Goldfarb (eds), *The Economics of Artificial Intelligence: An Agenda*, Chapter 21, University of Chicago Press, Chicago.
- Biau, G. and Scornet, E. (2016), A random forest guided tour, *TEST* 25, 197–227.
- Björklund, A., Jäntti, M., and Roemer, J. E. (2012), Equality of opportunity and the distribution of long-run income in Sweden, *Social Choice and Welfare* 39, 675–696.
- Blackburn, M. L. (2007), Estimating wage differentials without logarithms, *Labour Economics* 14, 73–98.
- Blau, F. D. and Kahn, L. M. (2017), The gender wage gap: extent, trends, and explanations, *Journal of Economic Literature* 55, 789–865.
- Blundell, J. and Risa, E. (2019), Income and family background: are we using the right models?, Working paper, available at <https://doi.org/10.2139/ssrn.3269576>.
- Bourguignon, F., Ferreira, F. H. G., and Menéndez, M. (2007), Inequality of opportunity in Brazil, *Review of Income and Wealth* 53, 585–618.
- Bratberg, E., Davis, J., Mazumder, B., Nybom, M., Schnitzlein, D. D., and Vaage, K. (2017), A comparison of intergenerational mobility curves in Germany, Norway, Sweden, and the US, *Scandinavian Journal of Economics* 119, 72–101.
- Breiman, L. (2001), Random forests, *Machine Learning* 45, 5–32.
- Breiman, L., Friedman, J., Stone, C. J., and Olshen, R. (1984), *Classification and Regression Trees*, Taylor & Francis, London.
- Brunori, P. and Neidhöfer, G. (2021), The evolution of inequality of opportunity in Germany: a machine learning approach, *Review of Income and Wealth* 67, 900–927.
- Brunori, P., Peragine, V., and Serlenga, L. (2019), Upward and downward bias when measuring inequality of opportunity, *Social Choice and Welfare* 52, 635–661.
- Brzezinski, M. (2020), The evolution of inequality of opportunity in Europe, *Applied Economics Letters* 27, 262–266.
- Cappelen, A. W., Hole, A. D., Sørensen, E. Ø., and Tungodden, B. (2007), The pluralism of fairness ideals: an experimental approach, *American Economic Review* 97 (3), 818–827.
- Carmichael, F., Darko, C. K., Ercolani, M. G., Ozgen, C., and Siebert, W. S. (2020), Evidence on intergenerational income transmission using complete Dutch population data, *Economics Letters* 189, 108996.
- Carranza, R. (2020), Upper and lower bound estimates of inequality of opportunity: a cross-national comparison for Europe, ECINEQ Working Paper Series 511.
- Checchi, D. and Peragine, V. (2010), Inequality of opportunity in Italy, *Journal of Economic Inequality* 8, 429–450.
- Checchi, D., Peragine, V., and Serlenga, L. (2016), Inequality of opportunity in Europe: is there a role for institutions?, in L. Cappellari, S. W. Polachek, and K. Tatsiramos (eds), *Inequality: Causes and Consequences*, Vol. 43, Chapter 1, Emerald Insight, Bingley, 1–44.
- Chetty, R., Hendren, N., Kline, P., and Saez, E. (2014a), Where is the land of opportunity? The geography of intergenerational mobility in the United States, *Quarterly Journal of Economics* 129, 1553–1623.
- Chetty, R., Hendren, N., Kline, P., Saez, E., and Turner, N. (2014b), Is the United States still a land of opportunity? Recent trends in intergenerational mobility, *American Economic Review* 104 (5), 141–147.
- Chetty, R., Hendren, N., Lin, F., Majerovitz, J., and Scuderi, B. (2016), Childhood environment and gender gaps in adulthood, *American Economic Review* 106 (5), 282–288.
- Corak, M. (2013), Income inequality, equality of opportunity, and intergenerational mobility, *Journal of Economic Perspectives* 27 (3), 79–102.
- Cowell, F. A. (2016), Inequality and poverty measures, in M. D. Adler and M. Fleurbaey (eds), *Oxford Handbook of Well-Being and Public Policy*, Chapter 4, Oxford University Press, Oxford, 82–125.

- Cowell, F. A. and Victoria-Feser, M.-P. (1996), Robustness properties of inequality measures, *Econometrica* 64, 77–101.
- Dahl, G. B. and Lochner, L. (2012), The impact of family income on child achievement: evidence from the Earned Income Tax Credit, *American Economic Review* 102 (5), 1927–1956.
- Ferreira, F. H. G. and Gignoux, J. (2011), The measurement of inequality of opportunity: theory and an application to Latin America, *Review of Income and Wealth* 57, 622–657.
- Friedberg, R., Tibshirani, J., Athey, S., and Wager, S. (2020), Local linear forests, *Journal of Computational and Graphical Statistics* 30, 503–517.
- Friedman, J., Hastie, T., and Tibshirani, R. (2009), *The Elements of Statistical Learning*, Springer, New York.
- Hothorn, T., Hornik, K., and Zeileis, A. (2006), Unbiased recursive partitioning: a conditional inference framework, *Journal of Computational and Graphical Statistics* 15, 651–674.
- Hufe, P., Peichl, A., Roemer, J. E., and Ungerer, M. (2017), Inequality of income acquisition: the role of childhood circumstances, *Social Choice and Welfare* 143, 499–544.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013), *An Introduction to Statistical Learning with Applications in R*, Springer, New York.
- Kanbur, R. and Snell, A. (2019), Inequality measures as tests of fairness, *Economic Journal* 129, 2216–2239.
- Kreisman, D. and Rangel, M. A. (2015), On the blurring of the color line: wages and employment for black males of different skin tones, *Review of Economics and Statistics* 97, 1–13.
- Lanza, S. T., Tan, X., and Bray, B. C. (2013), Latent class analysis with distal outcomes: a flexible model-based approach, *Structural Equation Modeling: A Multidisciplinary Journal* 20, 1–26.
- Lefranc, A., Pistolesi, N., and Trannoy, A. (2009), Equality of opportunity and luck: definitions and testable conditions, with an application to income in France, *Journal of Public Economics* 93, 1189–1207.
- Li Donni, P., Rodríguez, J. G., and Rosa Dias, P. (2015), Empirical definition of social types in the analysis of inequality of opportunity: a latent classes approach, *Social Choice and Welfare* 44, 673–701.
- Milanovic, B. (2015), Global inequality of opportunity: how much of our income is determined by where we live?, *Review of Economics and Statistics* 97, 452–460.
- Morgan, J. N. and Sonquist, J. A. (1963), Problems in the analysis of survey data, and a proposal, *Journal of the American Statistical Association* 58, 415–434.
- Mullainathan, S. and Spiess, J. (2017), Machine learning: an applied econometric approach, *Journal of Economic Perspectives* 31 (2), 87–106.
- Palomino, J. C., Marrero, G. A., and Rodríguez, J. G. (2019), Channels of inequality of opportunity: the role of education and occupation in Europe, *Social Indicators Research* 143, 1045–1074.
- Pistolesi, N. (2009), Inequality of opportunity in the land of opportunities, 1968–2001, *Journal of Economic Inequality* 7, 411–433.
- Ramos, X. and Van de gaer, D. (2016), Empirical approaches to inequality of opportunity: principles, measures, and evidence, *Journal of Economic Surveys* 30, 855–883.
- Roemer, J. E. (1998), *Equality of Opportunity*, Harvard University Press, Cambridge, MA.
- Roemer, J. E. and Trannoy, A. (2015), Equality of opportunity, in A. B. Atkinson and F. Bourguignon (eds), *Handbook of Income Distribution*, Vol. 2A, Chapter 4, Elsevier, Amsterdam, 217–300.
- Schafer, J. L. (1999), Multiple imputation: a primer, *Statistical Methods in Medical Research* 8, 3–15.
- Stuhler, J. (2018), A review of intergenerational mobility and its drivers, Joint Research Centre Technical Report, Publications Office of the European Union, Luxembourg.

- Suárez Álvarez, A. and López Menéndez, A. J. (2021), Dynamics of inequality and opportunities within European countries, *Bulletin of Economic Research* 73, 555–579.
- Van de gaer, D. (1993), Equality of opportunity and investment in human capital, PhD thesis, University of Leuven.
- Varian, H. R. (2014), Big data: new tricks for econometrics, *Journal of Economic Perspectives* 28 (2), 3–27.

First version submitted June 2020;
final version received February 2023.