

A Service of

ZBW

Leibniz-Informationszentrum Wirtschaft Leibniz Information Centre for Economics

Brunori, Paolo; Hufe, Paul; Mahler, Daniel Gerszon

Working Paper The Roots of Inequality: Estimating Inequality of Opportunity from Regression Trees and Forests

IZA Discussion Papers, No. 14689

Provided in Cooperation with: IZA – Institute of Labor Economics

Suggested Citation: Brunori, Paolo; Hufe, Paul; Mahler, Daniel Gerszon (2021) : The Roots of Inequality: Estimating Inequality of Opportunity from Regression Trees and Forests, IZA Discussion Papers, No. 14689, Institute of Labor Economics (IZA), Bonn

This Version is available at: https://hdl.handle.net/10419/245740

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.



WWW.ECONSTOR.EU



Initiated by Deutsche Post Foundation

DISCUSSION PAPER SERIES

IZA DP No. 14689

The Roots of Inequality: Estimating Inequality of Opportunity from Regression Trees and Forests

Paolo Brunori Paul Hufe Daniel Mahler

AUGUST 2021



Initiated by Deutsche Post Foundation

DISCUSSION PAPER SERIES

IZA DP No. 14689

The Roots of Inequality: Estimating Inequality of Opportunity from Regression Trees and Forests

Paolo Brunori

London School of Economics

Paul Hufe ifo Munich, LMU Munich and IZA

Daniel Mahler World Bank

AUGUST 2021

Any opinions expressed in this paper are those of the author(s) and not those of IZA. Research published in this series may include views on policy, but IZA takes no institutional policy positions. The IZA research network is committed to the IZA Guiding Principles of Research Integrity.

The IZA Institute of Labor Economics is an independent economic research institute that conducts research in labor economics and offers evidence-based policy advice on labor market issues. Supported by the Deutsche Post Foundation, IZA runs the world's largest network of economists, whose research aims to provide answers to the global labor market challenges of our time. Our key objective is to build bridges between academic research, policymakers and society.

IZA Discussion Papers often represent preliminary work and are circulated to encourage discussion. Citation of such a paper should account for its provisional character. A revised version may be available directly from the author.

IZA – Institute of Labor Economics					
Schaumburg-Lippe-Straße 5–9 53113 Bonn, Germany	Phone: +49-228-3894-0 Email: publications@iza.org	www.iza.org			

ABSTRACT

The Roots of Inequality: Estimating Inequality of Opportunity from Regression Trees and Forests^{*}

In this paper we propose the use of machine learning methods to estimate inequality of opportunity. We illustrate how our proposed methods—conditional inference regression trees and forests—represent a substantial improvement over existing estimation approaches. First, they reduce the risk of ad-hoc model selection. Second, they establish estimation models by trading 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 may 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. Our results illustrate the practical importance of leveraging machine learning algorithms to avoid giving misleading information about the level of inequality of opportunity in different societies to policymakers and the general public.

JEL Classification:D31, D63, C38Keywords:equality of opportunity, machine learning, random forests

Corresponding author:

Paolo Brunori International Inequalities Institute London School of Economics Houghton Street London, WC2A 2AE United Kingdom E-mail: paolo.brunori@unifi.it

^{*} We thank Chiara Binelli, Marc Fleurbaey, 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.

1 INTRODUCTION

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

In this paper we propose the use of machine learning methods to overcome the issue of adhoc 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 inequality of opportunity estimates that are less prone to upward or downward bias.

The empirical literature on the measurement of unequal opportunities has been flourishing since John Roemer's (1998) seminal contribution, *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 *effort*, and those factors are inequitable and call for compensation.¹ Grounded on this distinction, inequality of opportunity measures 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; and the more inequality of opportunity there is.

¹The distinction between circumstances and efforts underpins many prominent literature branches in economics 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).

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 increasing availability 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 (Chetty et al., 2016; Dahl and Lochner, 2012). 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 may 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) and Breiman (2001), they belong to a set of machine learning methods that is increasingly integrated into the statistical toolkit of economists (Varian, 2014; Mullainathan

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). In empirical implementations typical circumstances include biological sex, socioeconomic background, and race.

and Spiess, 2017; Athey, 2018). Trees and forests obtain predictions by drawing on a clearcut algorithm that imposes only minimal assumptions about which and how circumstances interact in shaping individual opportunities. 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 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 at par in terms of their opportunity characteristics. Hence, standard estimation approaches may yield misleading information about the level of inequality of opportunity in different societies to policymakers and the general public alike.

The remainder of this paper is organized as follows: section 2 gives a brief introduction to current empirical approaches in the literature on inequality of opportunity. Section 3 introduces conditional inference regression trees and forests, and illustrates how to use them in the context of inequality of opportunity estimations. An empirical illustration based on simulated data and the EU Survey of Income and Living Conditions is contained in section 4. In this section we also highlight the particular advantages of tree- and forest-based estimation methods by comparing them to the prevalent estimation approaches in the literature. Section 5 concludes the paper.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

2 EMPIRICAL APPROACHES TO EQUALITY OF OPPORTUNITY

Theoretical Set-up and Notation. Consider a population $\mathcal{N} = \{1, ..., N\}$ and an associated vector of non-negative incomes $y = (y_1, ..., y_N)$. 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 follows:

$$y = d(\omega, \theta). \tag{1}$$

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

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 a recent overview. In this work we exclusively focus on *ex-ante utilitarian* measures of inequality of opportunity (Van de gaer, 1993; Checchi and Peragine, 2010). They 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} :

$$y^{C} = (y_{1}^{C}, ..., y_{i}^{C}, ..., y_{N}^{C}) = (\mathbb{E}[y_{1}|\omega_{1}], ..., \mathbb{E}[y_{i}|\omega_{i}], ..., \mathbb{E}[y_{N}|\omega_{N}]).$$
(2)

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

⁴The use of ML methods is not restricted to ex-ante utilitarian formulations but can be easily extended to alternative measures of inequality of opportunity.

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

- 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() 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() 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(). While we present our main results in terms of the Gini index, we provide robustness checks based on other inequality indexes in Supplementary Material S.2.

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 follows:

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

 $\mathbb{E}(y|\omega)$ captures unfair variation due to observed circumstances. The iid 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.

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, $\beta = \frac{E(y_{ic}|y_{ip})}{y_{ip}}$, where y_{ip} represents parental income as the sole circumstance. Hence, the functional applied to the distribution of conditional expectations can be written as $I() = \frac{1}{y_{ip}}$. 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} .

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 that faces circumstances ω_i ? The precise form of f() is a priori unknown. In the vast majority of empirical applications, researchers address this lack of 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:

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

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

$$\hat{y}_i^{\mathsf{C}} = \mu_{m(i)} = \frac{1}{N_m} \sum_{j=1}^{N_m} y_j, \ \forall j \in t_m, \ \forall t_m \in \mathcal{T}.$$
(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 interdependencies between circumstances and other non-linearities. To alleviate this shortcoming the researcher may integrate interaction terms and higher order polynomials into equation (4). However, such extensions remain at her discretion. Reversely, the non-parametric approach does not restrict the interdependent impact of circumstances. However, if the data is rich enough in information on circumstances,

⁶The logarithmic transformation is not innocuous as the marginal impact of circumstances on incomes may 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.

the researcher may be forced to reduce the observed circumstance vector to obtain statistically meaningful estimates of the relevant parameters.⁷ The necessary process of restricting the circumstance vector again remains at the researcher's discretion.

The previous discussion illustrates that common approaches leave the researcher to her 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: the researcher 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 will discuss how regression trees and forests address the outlined shortcomings of existing estimation approaches.

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

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 will outline in the following, they can be straightforwardly applied to inequality of opportunity estimations and solve the issue of model selection.

First, we will 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 will introduce 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 $x = (x^1, ..., x^k)$. They use the sample $S = \{(y_i, x_i)\}_{i=1}^S$ to divide the population into non-overlapping groups, $G = \{g_1, ..., g_m, ..., g_M\}$, 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 $y = (y_1, ..., y_i, ..., y_N)$ one obtains a vector of predicted values $\hat{y} = (\hat{f}(x_1), ..., \hat{f}(x_i), ..., \hat{f}(x_N))$, where

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

The mapping from regression trees to equality of opportunity estimation is straightforward. If

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

the input variables $x = (x^1, ..., x^k)$ 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.

Tree Construction. Regression trees partition the sample into *M* types by *recursive binary splitting*. Recursive binary splitting 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 $\omega_i^p \ge \tilde{\omega}^p$, 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).





Note: 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} .

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 4-step procedure:

- 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} .
 - ⇒ 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.

$$\omega^* = \operatorname*{argmin}_{\omega^p} \{ p^{\omega^p}_{adj.} : \omega^p \in \Omega, \ p = 1, ..., P \}.$$

- \Rightarrow If $p_{adi.}^{\omega^*} > \alpha^*$: 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, $p^{\omega_s^*}$.
 - ⇒ Split the sample based on ω^* , by choosing the splitting point *s* that yields the lowest *p*-value, i.e. $\tilde{\omega}^* = \underset{\omega^*_s}{\operatorname{argmin}} \{ p^{\omega^*_s} : \omega^*_s \in \Omega \}.$
- 4. Repeat steps 1.–3. for each of the resulting subsamples.

In words, conditional inference 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 Appendix Figure A.1.

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 mimicks 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.

3.2 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 data generating process that imposes interactions while ruling out the linear influence of circumstances. Third, trees make inefficient use of data since some of the circumstances $\omega^p \in \Omega$ are not used for the construction of the tree. However, circumstances may 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).

Forest Construction. Random forests create many trees and average over all of these when making predictions. Trees are constructed according to the same 4-step procedure outlined in the previous subsection. However, two tweaks are made. First, given the sample $S = \{(y_i, \omega_i)\}_{i=1}^{S}$ each tree is estimated on a random subsample $S' \subset S$. In our application, we randomly select half 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, using 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:

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

Equation (8) illustrates that individual predictions are a function of α^* —the significance level governing the implementation of splits, \bar{P}^* —the number of circumstances to be considered at each splitting point, and B^* —the number of subsamples drawn from the data. 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 Appendix A.

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

- 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 $\hat{f}^{-H}()$.
- 3. Calculate the mean squared error in the test set: $MSE^{Test} = \frac{1}{N^{H}} \sum_{i \in H} [y_i - \hat{f}^{-H}(\omega_i)]^2.$

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 $\hat{f}()$ 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 indexes are presented in Supplementary Material S.2.

4.1 Data

We base our empirical illustration on the 2011 wave of the European Union Statistics on Income and Living Conditions (EU-SILC). EU-SILC 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 European Union. 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 is collected by national statistical agencies following common variable definitions and data collection procedures. We use the 2011 wave

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 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 Appendix **B**.

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.

since 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 respondent's characteristics at birth and their living conditions during childhood. Descriptive statistics of circumstance variables are reported in Supplementary Material S.1.

TABLE 1 – List of Circumstances

- 1. Respondent's sex:
 - Male
 - Female
- 2. Respondent's country of birth:
 - Respondent's 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/citizen of the respondent's present country of residence
 - Born/citizen of another EU-27 country
 - Born/citizen of another European country
 - Born/citizen of a country outside Europe
- 8. Father's/mother's education (based on the International Standard Classification of Education 1997 [ISCED-97])*:
 - Unknown father/mother
 - Illiterate
 - Low (0-2 ISCED-97)

- 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 the International Standard Classification of Occupations, published by the International Labour Office [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 occupation (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
- Medium (3-4 ISCED-97) - Not owned

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

Note: 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 11 is missing for Finland.

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 59. To assure the representativeness of our country samples we use individual cross-sectional weights throughout the analysis.

Table 2 shows considerable heterogeneity in income distributions across Europe. While the average households in Norway and Switzerland obtained incomes above \leq 40,000 in 2010, the average household income in Bulgaria, Romania and Lithuania did not exceed the \leq 5,000 mark. Lowest inequality prevails in Norway, Sweden and Iceland, all of which have Gini coefficients below 0.220. 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 (5 categories each), father's occupation (11 categories), area of birth (3 categories), and tenancy status of the household at age 14 (2 categories). The prediction model includes 22 parameters.

Second, we draw on the non-parametric approach as proposed by Checchi and Peragine (2010).

		Equivalized Disposable Household Income in \in			
Country	Ν	μ	σ	Gini	
Austria	6,220	25,451	13,971	0.268	
Belgium	6,011	23,291	10,948	0.249	
Bulgaria	7,154	3,714	2,491	0.333	
Croatia	6,969	6,627	3,819	0.306	
Cyprus	4,589	21,058	11,454	0.279	
Czech Republic	8,711	9,006	4,320	0.250	
Denmark	5,897	32,027	13,836	0.232	
Estonia	5,338	6,922	3,912	0.330	
Finland	9,743	27,517	13,891	0.246	
France	11,078	24,299	14,583	0.288	
Germany	12,683	22,221	12,273	0.276	
Greece	6,184	13,184	8,651	0.334	
Hungary	13,330	5,327	2,863	0.276	
Iceland	3,684	22,190	9,232	0.210	
Ireland	4,318	24,867	14,307	0.296	
Italy	21,070	18,786	11,730	0.309	
Latvia	6,423	5,334	3,618	0.363	
Lithuania	5,403	4,774	3,150	0.344	
Luxembourg	6,765	37,911	19,977	0.271	
Malta	4,701	13,006	6,747	0.277	
Netherlands	11,411	25,210	11,414	0.235	
Norway	5,026	43,260	16,971	0.202	
Poland	15,545	6,103	3,690	0.316	
Portugal	5,899	10,781	7,296	0.334	
Romania	7,867	2,562	1,646	0.337	
Slovakia	6,779	7,304	3,416	0.257	
Slovenia	13,183	13,772	5,994	0.225	
Spain	15,481	17,088	10,597	0.329	
Sweden	6,599	26,346	10,700	0.215	
Switzerland	7,583	42,208	24,486	0.279	
United Kingdom	7,391	25,936	16,815	0.320	

TABLE 2 – Summary Statistics

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

Note: 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.

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 (5 categories), fathers' occupation (4 categories), and migration status (2 categories).¹¹ 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 full set of observable circumstances from which the algorithm infers the appropriate number of unobserved types in the data by minimizing BIC.

Do these specification choices serve for a fair assessment of these benchmark methods? As

¹¹To 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.

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. We therefore 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 BIC to select the number of parameters to be estimated. We conclude that our specification choices for all three benchmark methods have clear precedence in the existing literature.

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.

We impose three DGP that are summarized in Table 3. The *parametric DGP* and *non-parametric DGP* correspond to the estimation models outlined in section 4.2. They present a challenging test for data-driven estimation methods since 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

	Parametric	Non-Parametric	Mixed
Outcome	ln(y)	У	ln(y)
Parameters	22	40	18
Circumstances	Education Father Education Mother Occupation Father Birth Area Tenancy Status	Education Parents Occupation Father Mig. Background	Education Parents Occupation Father Mig. Background Tenancy Status
Non-Linearity	None	Full Interaction	All Circ. w/ Mig. Background (2 Levels)
e	$\mathcal{N}(0, 2000)$	$\mathcal{N}(0, 2000)$	$\mathcal{N}(0, 2000)$

TABLE 3 – Summary of Data Generating Processes

DGP. This is a more realistic scenario since 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 = 200,754). In turn, we retain the predictions from these estimations and add a disturbance term with $\mathcal{N}(0,2000)$.¹² Thus, we obtain three variables that define the distribution of income for the purpose of this simulation.

Next we specify five sample sizes that broadly cover the range of effective country sample sizes observed in EU-SILC (see Table C.1): $N \in \{1,000;2,000;4,000;8,000;16,000\}$. For each sample size, we draw one test set of size $N^H = 1/3N$ and 50 training sets of size $N^{-H} = 2/3N$. 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: i) expected MSE^{Test} to assess out-of-sample prediction accuracy (James et al., 2013), and ii) the expected absolute difference between inequality of opportunity estimates and the true level of inequality of opportunity.

Figure 2 displays the results. In its lower part, each panel describes expected MSE^{Test} per combination of DGP, estimation approach and sample size. Since we know the true DGP we can decompose MSE^{Test} into variance and expected bias.¹³ In its upper part, each panel describes the corresponding absolute bias in inequality of opportunity estimates on an inverse scale. The

¹²We choose a variance term small enough such that $y \in \mathbb{R}_{++}$.

¹³See also Appendix **B** for an illustration of the variance-bias decomposition. The irreducible error term is uninformative for differences in MSE^{Test} because $Var(\epsilon) = 2,000^2$ is constant across specifications. Therefore, we only present evidence on the variance and the bias component of MSE^{Test}.

FIGURE 2 – Simulation Results





Data: EU-SILC 2011 cross-sectional (rev.5, June 2015). **Note:** This figure shows expected MSE^{Test} and expected bias in inequality of opportunity estimates from the simulation exercise. Each row corresponds to one data generating process (see Table 3). Each column corresponds to one estimation method (see sections 3.1, 3.2, 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 data generating process.

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 statistical theory. First, if the true DGP is known, expected bias is zero and MSE^{Test} is driven by its variance component only. Second, with increasing *N* the variance component of MSE^{Test} decreases for each combination of DGP and estimation approach. Third, 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, forests tend to have lower variance than trees—in our simulation this is true in 80% of all cases.

In terms of substantive results, it is clear that trees and forests dominate all other estimation approaches in terms of expected MSE^{Test}. This result holds both in comparison to fixed estimation approaches that do not invoke the true DGP (parametric, non-parametric) and in comparison to LCA as an alternative data-driven estimation approach. However, 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 furthermore highlight the close correspondence between MSE^{Test} and expected bias in inequality of opportunity estimates: the higher MSE^{Test}, the stronger inequality of opportunity estimates diverge from the ground truth.

In summary, the simulation results provide a strong case for 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.¹⁴

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 to 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}.

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

¹⁴We note that our 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 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.

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 to 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 14.2 (9.2%) (Figure 3, Panel [a]). With average shortfalls of 4.0 (2.5%) and 5.2 (2.9%), prediction errors are less pronounced for non-parametric (Figure 3, Panel [b]) and latent class models (Figure 3, Panel [c]). These averages, however, mask considerable heterogeneity. For example, the relative test error of parametric estimates for Slovenia, Slovakia and Czech Republic is below 4.0. To the contrary, the relative test error of parametric estimates for the UK and Luxembourg are more than eight times as large.

Conditional inference trees are closest to the test error rate of forests ($MSE^{Test} = 3.7 [2.3\%]$). Yet, they also fall short of the performance of forests due to higher variance, imposing non-linearity, and omission of less relevant circumstances (see section 3.2).

We conclude that among all considered methods, conditional inference forests deliver the highest out-of-sample prediction accuracy. 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.

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, since 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.¹⁵

¹⁵It 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. Methods that are upward biased in this comparison may potentially be closer to the ground truth than our reference estimate. Such statements,



FIGURE 3 – Comparison of MSE^{Test} by Method

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

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 indicates unweighted cross-country averages. Point estimates and associated standard errors are listed in Table D.1.

Panel (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

however, are purely speculative and cannot be falsified until better data becomes available (see also Appendix B 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

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

relative to forests which is the method providing the lowest expected bias in comparison to the true DGP. For 27 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 as high-opportunity societies. For example, forests classify Iceland and the Netherlands as the societies offering the highest equality of opportunity in Europe. To the

Note: 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 \hat{y}^{C} . Point estimates and associated standard errors are listed in Table D.2.

contrary, the parametric estimate would rank them at similar levels as France and Germany.

Panel (b) shows estimates from the non-parametric approach. The overall pattern resembles the parametric approach but on a more attenuated level. For 19 out of 31 countries the non-parametric estimate exceeds its forest-based analogue. These overstatements are again clustered among countries that are typically considered as high-opportunity societies. To the contrary, in the upper part of the equal-opportunity ranking the non-parametric approach tends to replicate the results from forests reasonably well. As evidenced by the simulation in section 4.3 such resemblance should be interpreted as a luck of the draw rather than a property inherent to the estimation approach. Under alternative type partitions, estimates from the non-parametric approach may diverge much more strongly than under the partition adopted in this work.

Panel (c) shows estimates from the latent class approach. The overall pattern is more heterogeneous than for the previous approaches. While overstatements prevail in countries that are typically considered as high-opportunity societies, there are 22 out of 31 countries for which the latent class estimate falls short of the forest estimate. These countries are clustered in the center and lower part of the the equal-opportunity ranking. For example, forests classify Germany and the UK at an intermediate position among all countries in Europe. To the contrary, the latent class estimate would elevate them into the top group next to the Nordic countries.

Panel (d) shows that trees and forests tend to produce similar results. The correlation between point estimates is high (0.99). In contrast to all other approaches there is no general tendency to over- or underestimate inequality of opportunity relative to forests.

Finally, all benchmark estimation approaches underestimate inequality of opportunity in the joint European sample. This finding emphasizes that our methodological conclusions about benchmark estimation approaches are contingent on their particular implementation in a given data environment. For example, consider the parametric approach that tends to overfit the data in the smaller country samples. In the large joint sample, upward bias due to noisy parameter estimates vanishes, while downward bias due to model mis-specification emerges as the dominant force driving differences to forests. In view of such contingency, one may ask to what extent our conclusions on inequality of opportunity in Europe are driven by differences and

sample sizes across methods and countries. Therefore, we assess the robustness of our results to sample size differences in the following.

Robustness to Differences in Sample Size. Effective sample sizes differ by estimation method and country (Table C.1). 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% in 6 countries of our sample (Denmark, Iceland, Netherlands, Norway, Slovenia, 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 = 3,684). Therefore, we perform two robustness analyses.

First, we recompute inequality of opportunity after completing missing data through multiple imputation (Schafer, 1999).¹⁶ As a consequence, we can compare inequality of opportunity estimates across methods on the same effective sample size per country. Figure C.1 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. To the contrary, the patterns for trees and forests remain unaffected since they handle missing values by default through surrogate splits. The general pattern of our methodological comparison remains unaffected.

Second, we recompute inequality of opportunity 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 C.2 shows that point estimates and country rankings differ strongly for all benchmark methods (parametric, non-parametric, LCA). To the contrary, point estimates and country rankings of trees and forests are unaffected by harmonization in sample sizes across countries. This result bolsters confidence that opportunity rankings of trees and forests are not an artifact of cross-country variation in

¹⁶List-wise deletion yields unbiased parameter estimates if data is missing completely at random (MCAR). Multiple imputation weakens this assumption by assuming that data is missing at random (MAR), i.e. that missing data is random conditional on observed variables.

sample sizes.

Comparison to 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 may be misleading in two related dimensions. First, they may 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 may turn to the wrong country examples. In the following we will 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 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, researcher discretion in model selection would be irrelevant for conclusions about inequality of opportunity in Europe. Second, we assess whether 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 since 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, 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,¹⁷ ii) existing consensus estimates of the IGE,¹⁸ and iii) inequality of opportunity estimates from our preferred methods—regression trees and forests.

¹⁷We 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), Carranza (2020), and Hufe et al. (2018). Furthermore, we do not include Brzezinski (2020) since he derives estimates based on the methods proposed in this paper.

¹⁸We focus on IGE estimates based on actual data linkages across generations and 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.

Panel (a) of Table 4 suggests that existing literature on inequality of opportunity in Europe is not consistent. Rank correlations as low as 0.09 indicate strong heterogeneity in country rankings. This result is notable since all estimates were derived from the same underlying data source (EU-SILC), refer to a similar age group (approx. 25-60), and summarize counterfactual distributions \hat{y}^{C} by the same inequality metric (mean log deviation). Hence, discretionary choices with respect to model specifications may be a major force behind inconclusive evidence in the inequality of opportunity literature.¹⁹ To the contrary, application of our preferred estimation methods explicitly addresses this source of incoherence.

	Existing Studies		This Paper					
	Checchi et al. (2016)	Palomino et al. (2019)	Suarez et al. (2021)	Tree	Forest			
Panel (a): Equality of Opportunity (23 countries)								
Tree	•			1.000				
Forest				0.983	1.000			
Checchi et al. (2016)	1.000			0.388	0.380			
Palomino et al. (2019)	0.281	1.000		0.877	0.875			
Suarez et al. (2021)	0.090	0.855	1.000	0.738	0.762			
Panel (b): Intergenerational Elasticity (10 countries)								
Stuhler (2018) & Carmichael et al. (2020)	0.535	0.657	0.444	0.894	0.869			

TABLE 4 – Rank Correlations of Existing Studies

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

Note: 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, Palomino et al. (2019), Checchi et al. (2016), 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 .

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 moderately correlated with IGE rankings. However, various findings contradict comparative evidence on the IGE (Carmichael

¹⁹We acknowledge that differences differences in income concept definitions, i.e. individual vs. household income, may also contribute to the observed divergence.

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 at par with the Nordic countries. Checchi et al. (2016) find the Netherlands to be in the lower part of the opportunity ranking. To 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 data generating processes 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.

Next to these advantages, we acknowledge two potential drawbacks of our preferred methods for empirical research. First, (non-)parametric estimation approaches can be estimated by OLS—one of the workhorse estimation methods in economics and other social sciences. To the contrary, machine learning tools may 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.²⁰ Therefore, time-constrained researchers who are willing to settle for a fair first-order approximation of the least biased method may consider using trees instead of forests.

To be sure, 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; Kanbur and Snell, 2018; Pistolesi, 2009; Brunori and Neidhöfer, 2021).

²⁰The 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 a AMD Ryzen 7 4700U Processor (8 cores) and 16 GB RAM working memory.

References

- ALESINA, Alberto, Stefanie STANTCHEVA, and Edoardo TESO (2018). "Intergenerational Mobility and Preferences for Redistribution". *American Economic Review* 108 (2), pp. 521–554.
- ANDREOLI, Francesco and Alessio FUSCO (2019). "Robust cross-country analysis of inequality of opportunity". *Economics Letters* 182, pp. 86–89.
- ARNESON, Richard J. (2018). "Four Conceptions of equal opportunity". *Economic Journal* 128 (612), F152–F173.
- ATHEY, Susan (2018). "The Impact of Machine Learning on Economics". *The Economics of Artificial Intelligence: An Agenda*. Ed. by Ajay K. AGRAWAL, Joshua GANS, and Avi GOLDFARB. Chicago: University of Chicago Press. Chap. 21.
- BIAU, Gérard and Erwan SCORNET (2016). "A random forest guided tour". *Test* 25 (2), pp. 197–227.
- BJÖRKLUND, Anders, Markus JÄNTTI, and John E. ROEMER (2012). "Equality of opportunity and the distribution of long-run income in Sweden". *Social Choice and Welfare* 39 (2-3), pp. 675– 696.
- BLACKBURN, McKinley L. (2007). "Estimating wage differentials without logarithms". *Labour Economics* 14 (1), pp. 73–98.
- BLAU, Francine D. and Lawrence M. KAHN (2017). "The Gender Wage Gap: Extent, Trends, and Explanations". *Journal of Economic Literature* 55 (3), pp. 789–865.
- BLUNDELL, Jack and Erling RISA (2019). "Income and family background: Are we using the right models?" *mimeo*.
- BOURGUIGNON, François, Francisco H. G FERREIRA, and Marta MENÉNDEZ (2007). "Inequality of Opportunity in Brazil". *Review of Income and Wealth* 53 (4), pp. 585–618.
- BRATBERG, Espen, Jonathan DAVIS, Bhashkar MAZUMDER, Martin NYBOM, Daniel D. SCHNIT-ZLEIN, and Kjell VAAGE (2017). "A Comparison of Intergenerational Mobility Curves in Germany, Norway, Sweden, and the US". *Scandinavian Journal of Economics* 119 (1), pp. 72–101.
- BREIMAN, Leo (2001). "Random Forests". Machine Learning 45 (1), pp. 5–32.
- BREIMAN, Leo, Jerome FRIEDMAN, Charles J. STONE, and R.A. OLSHEN (1984). *Classification and Regression Trees*. Belmont: Taylor & Francis.

- BRUNORI, Paolo and Guido NEIDHÖFER (2021). "The Evolution of Inequality of Opportunity in Germany: A Machine Learning Approach". *Review of Income and Wealth*, Forthcoming.
- BRUNORI, Paolo, Vito PERAGINE, and Laura SERLENGA (2019). "Upward and downward bias when measuring inequality of opportunity". *Social Choice and Welfare* 52, pp. 635–661.
- BRZEZINSKI, Michal (2020). "The evolution of inequality of opportunity in Europe". *Applied Economics Letters* 27 (4), pp. 262–266.
- CAPPELEN, Alexander W., Astri Drange HOLE, Erik Ø. SØRENSEN, and Bertil TUNGODDEN (2007). "The Pluralism of Fairness Ideals: An Experimental Approach". *American Economic Review* 97 (3), pp. 818–827.
- CARMICHAEL, Fiona, Christian K. DARKO, Marco G. ERCOLANI, Ceren OZGEN, and W. Stanley SIEBERT (2020). "Evidence on intergenerational income transmission using complete Dutch population data". *Economics Letters* 189, p. 108996.
- CARRANZA, Rafael (2020). "Upper and lower bound estimates of inequality of opportunity: A cross-national comparison for Europe". ECINEQ Working Paper Series 511.
- CHECCHI, Daniele and Vito PERAGINE (2010). "Inequality of opportunity in Italy". *Journal of Economic Inequality* 8 (4), pp. 429–450.
- CHECCHI, Daniele, Vito PERAGINE, and Laura SERLENGA (2016). "Inequality of Opportunity in Europe: Is There a Role for Institutions?" *Inequality: Causes and Consequences*. Ed. by Lorenzo CAPPELLARI, Solomon W. POLACHEK, and Konstantinos TATSIRAMOS. Vol. 43. Emerald Insight. Chap. 1, pp. 1–44.
- CHETTY, Raj, Nathaniel HENDREN, Patrick KLINE, and Emmanuel SAEZ (2014a). "Where is the Land of Opportunity? The Geography of Intergenerational Mobility in the United States". *Quarterly Journal of Economics* 129 (4), pp. 1553–1623.
- CHETTY, Raj, Nathaniel HENDREN, Patrick KLINE, Emmanuel SAEZ, and Nicholas TURNER (2014b). "Is the United States Still a Land of Opportunity? Recent Trends in Intergenerational Mobility". *American Economic Review* 104 (5), pp. 141–47.
- CHETTY, Raj, Nathaniel HENDREN, Frina LIN, Jeremy MAJEROVITZ, and Benjamin SCUDERI (2016). "Childhood Environment and Gender Gaps in Adulthood". *American Economic Review* 106 (5), pp. 282–88.
- CORAK, Miles (2013). "Income Inequality, Equality of Opportunity, and Intergenerational Mobility". *Journal of Economic Perspectives* 27 (3), pp. 79–102.

- COWELL, Frank A. (2016). "Inequality and Poverty Measures". *Oxford Handbook of Well-Being and Public Policy*. Ed. by Matthew D. ADLER and Mark FLEURBAEY. Oxford: Oxford University Press. Chap. 4, pp. 82–125.
- COWELL, Frank A. and Maria-Pia VICTORIA-FESER (1996). "Robustness Properties of Inequality Measures". *Econometrica* 64 (1), pp. 77–101.
- DAHL, Gordon B. and Lance LOCHNER (2012). "The Impact of Family Income on Child Achievement: Evidence from the Earned Income Tax Credit". *American Economic Review* 102 (5), pp. 1927–56.
- FERREIRA, Francisco H. G and Jérémie GIGNOUX (2011). "The Measurement of Inequality of Opportunity: Theory and an Application to Latin America". *Review of Income and Wealth* 57 (4), pp. 622–657.
- FRIEDBERG, Rina, Julie TIBSHIRANI, Susan ATHEY, and Stefan WAGER (2020). "Local Linear Forests". *Journal of Computational and Graphical Statistics* 30 (2), pp. 503–517.
- FRIEDMAN, Jerome, Trevor HASTIE, and Robert TIBSHIRANI (2009). *The elements of statistical learning*. New York: Springer.
- HOTHORN, Torsten, Kurt HORNIK, and Achim ZEILEIS (2006). "Unbiased Recursive Partitioning: A Conditional Inference Framework". *Journal of Computational and Graphical Statistics* 15 (3), pp. 651–674.
- HUFE, Paul, Ravi KANBUR, and Andreas PEICHL (2018). "Measuring Unfair Inequality: Reconciling Equality of Opportunity and Freedom from Poverty". *CEPR Discussion Paper* 12989.
- HUFE, Paul, Andreas PEICHL, John E. ROEMER, and Martin UNGERER (2017). "Inequality of Income Acquisition: The Role of Childhood Circumstances". *Social Choice and Welfare* 143 (3-4), pp. 499–544.
- JAMES, Gareth, Daniela WITTEN, Trevor HASTIE, and Robert TIBSHIRANI (2013). *An Introduction to Statistical Learning with Applications in R.* New York: Springer.
- KANBUR, Ravi and Andy SNELL (2018). "Inequality Measures as Tests of Fairness". *Economic Journal* Forthcoming.
- KREISMAN, Daniel and Marcos A. RANGEL (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), pp. 1–13.
- LANZA, Stephanie T., Xianming TAN, and Bethany C. BRAY (2013). "Latent Class Analysis With Distal Outcomes: A Flexible Model-Based Approach". *Structural Equation Modelling* 20 (1), pp. 1–26.
- LEFRANC, Arnaud, Nicolas PISTOLESI, and Alain TRANNOY (2009). "Equality of opportunity and luck: Definitions and testable conditions, with an application to income in France". *Journal of Public Economics* 93 (11–12), pp. 1189–1207.
- LI DONNI, Paolo, Juan Gabriel RODRÍGUEZ, and Pedro ROSA DIAS (2015). "Empirical definition of social types in the analysis of inequality of opportunity: a latent classes approach". *Social Choice and Welfare* 44 (3), pp. 673–701.
- MORGAN, James N. and John A. SONQUIST (1963). "Problems in the Analysis of Survey Data, and a Proposal". *Journal of the American Statistical Association* 58 (302), pp. 415–434.
- MULLAINATHAN, Sendhil and Jann SPIESS (2017). "Machine Learning: An Applied Econometric Approach". *Journal of Economic Perspectives* 31 (2), pp. 87–106.
- PALOMINO, Juan C., Gustavo A. MARRERO, and Juan G. RODRÍGUEZ (2019). "Channels of inequality of opportunity: The role of education and occupation in Europe". *Social Indicators Research* 143, pp. 1045–1074.
- PISTOLESI, Nicolas (2009). "Inequality of opportunity in the land of opportunities, 1968–2001". *Journal of Economic Inequality* 7 (4), pp. 411–433.
- RAMOS, Xavier and Dirk VAN DE GAER (2016). "Empirical Approaches to Inequality of Opportunity: Principles, Measures, and Evidence". *Journal of Economic Surveys* 30 (5), pp. 855–883.
- ROEMER, John E. (1998). Equality of Opportunity. Cambridge: Harvard University Press.
- ROEMER, John E. and Alain TRANNOY (2015). "Equality of Opportunity". *Handbook of Income Distribution*. Ed. by Anthony B. ATKINSON and François BOURGUIGNON. Vol. 2A. Amsterdam: Elsevier. Chap. 4, pp. 217–300.
- SCHAFER, Joseph L (1999). "Multiple imputation: a primer". *Statistical Methods in Medical Research* 8 (1), pp. 3–15.
- STUHLER, Jan (2018). A Review of Intergenerational Mobility and its Drivers. Luxembourg: Publications Office of the European Union.
- SUÁREZ ÁLVAREZ, Ana and Ana Jesús LÓPEZ MENÉNDEZ (2021). "Dynamics of inequality and opportunities within European countries". *Bulletin of Economic Research*, Forthcoming.

- VAN DE GAER, Dirk (1993). "Equality of Opportunity and Investment in Human Capital". PhD thesis. University of Leuven.
- VARIAN, Hal R. (2014). "Big Data: New Tricks for Econometrics". *Journal of Economic Perspectives* 28 (2), pp. 3–27.

A EMPIRICAL CHOICES

Tuning of Trees. Alternatively to specifying α^* a priori, it can be chosen by *K*-fold cross-validation (CV), which—under some minimal assumptions—provides unbiased estimates of the out-of-sample MSE (Friedman et al., 2009). First, one splits the sample into *K* equal-sized folds. Second, one implements the conditional inference algorithm on the union of *K* – 1 folds for varying levels of α . This step makes it possible to compare the prediction from the *K* – 1 folds with the unused data points in the *k*th fold. Third, one calculates the out-of-sample MSE as a function of α :

$$MSE_k^{CV}(\alpha) = \frac{1}{N^k} \sum_{i \in k} (y_i^k - \hat{f}^{-k}(\omega_i; \alpha))^2, \ \omega_i \in \Omega, \ i \in \mathcal{N},$$
(9)

where $\hat{f}^{-k}()$ denotes the estimation function constructed while leaving out the *k*th fold. Fourth, this exercise is repeated for all *K* folds, so that $MSE^{CV}(\alpha) = \frac{1}{K} \sum_{k} MSE_{k}^{CV}(\alpha)$. Finally, one chooses α^* such that

$$\alpha^* = \operatorname*{argmin}_{\alpha} \{ \mathrm{MSE}^{\mathrm{CV}}(\alpha) : \alpha \in (0,1) \}.$$
(10)

Figure A.1 reveals that selecting α^* based on cross-validation or setting $\alpha^* = 0.05$ has little bearing on our results.

Tuning of Forests. The grid of parameters $(\alpha, \overline{P}, B)$ can be imposed a priori by the researcher or tuned to optimize the out-of-sample fit of the model. In our empirical illustration we proceed as follows. First, we fix B^* at a level at which the marginal gain of drawing an additional subsample in terms of out-of-sample prediction accuracy becomes negligible. Empirical tests show that this is the case with $B^* = 200$ for most countries in our sample (Figure A.2).

Second, we determine α^* and \bar{P}^* by minimizing the *out-of-bag* error. This entails the following three steps for a grid of values of α and \bar{P} :

1. Run a random forest with B^* subsamples, where \overline{P} circumstances are randomly chosen to be considered at each splitting point, and α is used as the critical value for the hypothesis tests.



FIGURE A.1 – Tuning Conditional Inference Trees

Data: EU-SILC 2011 cross-sectional (rev.5, June 2015). **Note:** This figure shows MSE^{Test} for different specifications of α^* relative to the baseline specification of $\alpha^* = 0.01$. Relative MSE^{Test} > 1 indicates worse fit than the baseline specification. Tuning is conducted by 5-fold CV. 95% confidence intervals are derived based on 200 bootstrapped re-samples of the test data using the normal approximation method.



FIGURE A.2 – Optimal Size of Forests

Data: EU-SILC 2011 cross-sectional (rev.5, June 2015). **Note:** This figure shows MSE^{OOB} for different specifications of B^* in the country sample of Germany. We fix $\bar{P}^* = 6$. The blue line is a non-parametric fit for MSE^{OOB} estimates. Gray shades indicate the 95% confidence interval.

- 2. Calculate the average predicted value of observation *i* using each of the prediction functions estimated in the subsamples $\mathcal{B}_{-i} := \{\mathcal{S}' \subset \mathcal{S} : \mathcal{S}' \cap \{(y_i, \omega_i)\} = \emptyset\}$ (the so called *bags*) in which *i* does not enter: $\hat{f}^{OOB}(\omega_i; \alpha, \bar{P}) = \frac{1}{N_{\mathcal{B}_{-i}}} \sum_{\mathcal{S}' \in \mathcal{B}_{-i}} \hat{f}^{\mathcal{S}'}(\omega_i; \alpha, \bar{P})$.
- 3. Calculate the out-of-bag mean squared error:

$$MSE^{OOB}(\alpha, \bar{P}) = \frac{1}{N} \sum_{i} [y_i - \hat{f}^{OOB}(\omega_i; \alpha, \bar{P})]^2.$$

Finally, one chooses the combination of parameter values that delivers the lowest MSE^{OOB}:

$$(\alpha^*, \bar{P}^*) = \operatorname*{argmin}_{\alpha, \bar{P}} \{ \mathrm{MSE}^{\mathrm{OOB}} : (\alpha, \bar{P}) \in (0, 1) \times \bar{\mathbf{P}} \}.$$
(11)

B UPWARD BIAS, DOWNWARD BIAS AND THE MSE

A standard statistic to assess prediction accuracy is the mean squared error (MSE):

$$MSE = \mathbb{E}_{\mathcal{S}}[(y - \hat{f}(\omega))^2], \qquad (12)$$

where *y* is the observed outcome and $\hat{f}(\omega)$ the estimator of the conditional expectation $\mathbb{E}(y|\omega)$ in a random sample S. The MSE can be decomposed into (1) variance, (2) expected bias, (3) irreducible error term (Friedman et al., 2009):

$$MSE = Var(\hat{f}(\omega)) + \mathbb{E}_{\mathcal{S}}[f(\omega) - \hat{f}(\omega)]^2 + Var(\epsilon),$$
(13)

$$=\underbrace{\operatorname{Var}(f(\omega) - \hat{f}(\omega))}_{(1)} + \underbrace{(f(\omega) - \mathbb{E}_{\mathcal{S}}[\hat{f}(\omega)])^2}_{(2)} + \underbrace{\operatorname{Var}(\epsilon)}_{(3)}.$$
(14)

In the literature on statistical learning this is referred to as the variance-bias decomposition. All three components can be linked to upward and downward biases in inequality of opportunity estimates.

(1) The variance captures upward bias due to model mis-specification. To see this, note that we minimize (1) by imposing the following model specification $y = \hat{f}(\omega) + \epsilon = \beta_0 + \epsilon$, i.e. by assuming that individual outcomes are best predicted by the sample mean μ^{S} .²¹ As a consequence, (1) drops out and MSE is entirely captured by components (2) and (3):

$$MSE = Var(f(\omega) - \hat{f}(\omega)) + (f(\omega) - \mathbb{E}_{\mathcal{S}}[\hat{f}(\omega)])^{2} + Var(\epsilon)$$
$$= (f(\omega) - \mu)^{2} + Var(\epsilon).$$

Hence, *the variance-minimizing estimation model* cannot yield upward biased estimates of inequality of opportunity since it is restricted in a way that does not allow for any role of ω in the explanation of y. To the contrary, it will be downward biased. For any functional I() that satisfies the measurement criteria outlined in section 2, $I(\hat{y}^C) = 0$.

²¹For the sake of exposition, we additionally assume $\mu^{S} = \mu$. Obviously, this is a stark assumption. In reality, there will always be some variance in sample means as long as one does not capture the entire population.

(2) The expected bias captures downward bias due to model mis-specification. To see this, note that we minimize (2) by specifying a complex model that allows for all observable circumstances, their mutual interactions and non-linearities.²² As a consequence, (2) drops out and MSE is entirely captured by components (1) and (3):

$$\begin{split} \mathrm{MSE} &= \mathrm{Var}(f(\omega) - \hat{f}(\omega)) + (f(\omega) - \mathbb{E}_{\mathcal{S}}[\hat{f}(\omega)])^2 + \mathrm{Var}(\epsilon), \\ &= \mathrm{Var}(f(\omega) - \hat{f}(\omega)) + \mathrm{Var}(\epsilon). \end{split}$$

Hence, in expectation and within a given data environment *the bias-minimizing estimation model* cannot yield downward biased estimates of inequality of opportunity. To the contrary, it will be upward biased. The conditional expectations within a particular sample S is estimated with error and measurement error inflates the variance of \hat{y}^C in comparison to the underlying truth: $Var(\hat{y}^C) = Var(y^C) + Var(u)$. For any functional I() that satisfies the measurement criteria outlined in section 2, $I(\hat{y}^C) > I(y^C)$.

(3) The irreducible error term contains downward bias due to unobserved circumstance variables. To see this, assume we relax the assumption that we observe the full set of relevant circumstances. In this case, variation due to unobserved circumstances is captured in the irreducible error term (3). This part of downward bias prevails regardless of estimation method and decreases as more circumstance information becomes available. Therefore, minimizing the out-of-sample MSE corresponds to minimizing expected bias in inequality of opportunity estimates *conditional on a given data environment*.

²²For the sake of exposition, we additionally assume that we observe all relevant circumstances. Obviously, this is a stark assumption. In reality, non-observable circumstance information is a key reason for downward bias in inequality of opportunity estimates that prevails regardless of estimation methods.

C SENSITIVITY TO SAMPLE SIZE

		Benchmark Methods	Conditional Inference				
Country	Parametric	Non-Parametric	Latent Class	Tree	Forest		
Austria	6,060	6,107	5,961	6,220	6,220		
Belgium	5,289	5,439	4,412	6,011	6,011		
Bulgaria	6,104	6,212	5,728	7,154	7,154		
Croatia	5,997	6,159	5,329	6,969	6,969		
Cyprus	4,491	4,525	4,448	4,589	4,589		
Czech Republic	6,488	6,524	5,826	8,711	8,711		
Denmark	2,218	2,302	1,985	5,897	5,897		
Estonia	4,918	5,004	4,696	5,338	5,338		
Finland	3,080	3,209	2,336	9,743	9,743		
France	10,214	10,433	9,710	11,078	11,078		
Germany	10,964	11,149	9,104	12,683	12,683		
Greece	5,767	5,862	5,516	6,184	6,184		
Hungary	12,324	12,526	12,019	13,330	13,330		
Iceland	1,481	1,552	1,393	3,684	3,684		
Ireland	3,102	3,164	2,880	4,318	4,318		
Italy	20,284	20,803	20,012	21,070	21,070		
Latvia	6,142	6,192	4,894	6,423	6,423		
Lithuania	4,613	4,705	4,344	5,403	5,403		
Luxembourg	6,567	6,654	6,361	6,765	6,765		
Malta	4,082	4,117	3,915	4,701	4,701		
Netherlands	5,461	5,598	5,135	11,411	11,411		
Norway	2,355	2,456	2,221	5,026	5,026		
Poland	12,808	13,369	12,498	15,545	15,545		
Portugal	5,696	5,809	5,624	5,899	5,899		
Romania	5,834	6,145	4,989	7,867	7,867		
Slovakia	6,212	6,404	6,049	6,779	6,779		
Slovenia	4,696	4,749	4,629	13,183	13,183		
Spain	14,672	14,817	14,473	15,481	15,481		
Sweden	531	624	439	6,599	6,599		
Switzerland	6,482	6,766	4,673	7,583	7,583		
United Kingdom	5,847	5,922	5,604	7,391	7,391		

TABLE C.1 – Sample Size by Method

Data: EU-SILC 2011 cross-sectional (rev.5, June 2015). **Note:** This table shows effective sample sizes for inequality of opportunity estimations by estimation method.



FIGURE C.1 – Inequality of Opportunity Estimates: Robustness to Multiple Imputation



Note: This figure shows inequality of opportunity estimates from different estimation methods relative to forests. We impute missing circumstance information by multiple imputation such that sample sizes are constant across methods. For each country we make 10 imputations, estimate inequality of opportunity, and calculate the corresponding average. Inequality of opportunity is measured by the Gini coefficient of the counterfactual distribution \hat{y}^{C} .



FIGURE C.2 – Inequality of Opportunity Estimates: Robustness to Sample Size Reductions

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

Note: This figure shows changes in inequality of opportunity estimates when reducing estimation samples to the smallest methodspecific sample size. For each country-method cell we make 10 random draws from the full country sample, estimate inequality of opportunity, and calculate the corresponding average. Inequality of opportunity is measured by the Gini coefficient of the counterfactual distribution \hat{y}^{C} .

D POINT ESTIMATES AND STANDARD ERRORS

	1	Benchmark Methods	Conditional Inference					
Country	Parametric	Non-Parametric	Latent Class	Tree	Forest			
Austria	281.1 [18.1]	275.3 [16.3]	280.1 [16.5]	272.6 [16.2]	268.5 [16.2]			
Belgium	163.7 [11.7]	157.2 [9.6]	160.9 [9.5]	156.2 [10.3]	152.3 [10.0]			
Bulgaria	37.6 [3.2]	35.7 [2.8]	36.4 [2.8]	35.9 [2.8]	35.2 [2.8]			
Croatia	51.6 [2.8]	46.2 [2.4]	46.0 [2.4]	46.6 [2.4]	45.9 [2.4]			
Cyprus	264.8 [21.5]	254.0 [19.9]	259.9 [20.0]	253.9 [20.0]	249.6 [19.7]			
Czech Republic	49.0 [3.2]	48.5 [2.9]	47.9 [2.9]	48.4 [3.0]	47.4 [2.9]			
Denmark	219.8 [25.6]	194.8 [21.0]	194.3 [21.6]	193.8 [20.7]	189.5 [20.2]			
Estonia	50.0 [3.1]	46.5 [2.7]	47.4 [2.9]	47.0 [2.8]	45.8 [2.8]			
Finland	253.0 [22.5]	242.4 [20.0]	236.7 [20.3]	235.9 [19.9]	234.8 [19.6]			
France	295.7 [19.3]	285.4 [17.7]	287.0 [17.8]	283.1 [17.6]	279.4 [17.2]			
Germany	262.0 [11.9]	249.4 [10.2]	250.3 [10.2]	251.1 [10.5]	245.0 [10.4]			
Greece	153.3 [12.9]	132.1 [11.2]	134.6 [11.7]	131.8 [11.2]	128.1 [10.9]			
Hungary	34.3 [1.4]	33.0 [1.2]	33.0 [1.2]	33.0 [1.2]	32.2 [1.2]			
Iceland	139.5 [20.2]	118.4 [16.7]	119.2 [17.0]	118.3 [16.8]	115.5 [16.5]			
Ireland	259.0 [22.1]	243.5 [19.7]	252.2 [19.6]	250.4 [19.7]	234.8 [18.6]			
Italy	238.1 [6.8]	214.4 [5.9]	218.3 [6.0]	210.6 [6.0]	207.2 [5.9]			
Latvia	43.0 [2.6]	38.2 [2.0]	39.0 [2.1]	39.4 [2.1]	37.6 [2.0]			
Lithuania	40.8 [2.5]	37.6 [2.1]	36.8 [2.1]	37.6 [2.1]	37.0 [2.0]			
Luxembourg	415.7 [26.0]	396.9 [23.0]	402.1 [23.3]	396.1 [23.4]	379.2 [22.6]			
Malta	114.9 [7.7]	108.4 [6.6]	107.8 [6.8]	109.3 [6.5]	106.1 [6.5]			
Netherlands	190.2 [12.2]	184.1 [11.2]	180.4 [11.1]	180.4 [11.0]	179.5 [11.1]			
Norway	205.6 [17.4]	198.4 [15.6]	197.7 [16.0]	196.3 [17.0]	197.6 [16.1]			
Poland	69.4 [3.3]	65.5 [2.9]	65.9 [3.0]	65.7 [2.9]	64.5 [2.9]			
Portugal	109.2 [8.3]	107.3 [7.3]	106.4 [7.4]	107.0 [7.2]	101.3 [6.9]			
Romania	14.3 [1.0]	14.3 [0.9]	13.6 [0.9]	13.5 [0.8]	13.2 [0.9]			
Slovakia	43.7 [2.4]	41.5 [2.1]	41.5 [2.1]	41.6 [2.3]	41.3 [2.2]			
Slovenia	83.1 [5.2]	80.9 [4.8]	82.0 [5.0]	80.9 [4.8]	80.5 [4.7]			
Spain	200.1 [7.4]	182.2 [6.3]	190.2 [6.4]	181.2 [6.3]	178.1 [6.3]			
Sweden	170.2 [29.6]	151.2 [24.4]	152.4 [27.2]	152.0 [25.8]	149.1 [25.5]			
Switzerland	403.6 [28.6]	392.3 [26.2]	393.0 [26.9]	397.7 [27.0]	382.4 [26.0]			
United Kingdom	423.7 [38.5]	386.1 [33.7]	386.2 [33.8]	388.7 [33.4]	380.0 [32.9]			

TABLE D.1 – MSE^{Test} Estimates

Data: EU-SILC 2011 cross-sectional (rev.5, June 2015). **Note:** This tables shows MSE^{Test} for different estimation methods. For all methods, we multiply MSE^{Test} by 1×10^{-6} . Standard errors are derived from 200 bootstrapped re-samples of the test data.

		Benchmark Methods	Conditional Inference						
Country	Parametric	Non-Parametric	Latent Class	Tree	Forest				
Austria	0.088 [0.004]	0.075 [0.004]	0.080 [0.005]	0.087 [0.005]	0.088 [0.004]				
Belgium	0.108 [0.006]	0.087 [0.004]	0.053 [0.009]	0.087 [0.005]	0.091 [0.003]				
Bulgaria	0.152 [0.005]	0.136 [0.005]	0.115 [0.005]	0.136 [0.005]	0.134 [0.005]				
Croatia	0.125 [0.007]	0.088 [0.005]	0.076 [0.006]	0.082 [0.006]	0.076 [0.005]				
Cyprus	0.093 [0.004]	0.083 [0.005]	0.074 [0.008]	0.080 [0.007]	0.080 [0.008]				
Czech Republic	0.066 [0.004]	0.066 [0.004]	0.060 [0.005]	0.057 [0.005]	0.051 [0.004]				
Denmark	0.070 [0.013]	0.041 [0.005]	0.029 [0.005]	0.021 [0.005]	0.020 [0.003]				
Estonia	0.108 [0.006]	0.102 [0.006]	0.074 [0.008]	0.097 [0.007]	0.101 [0.006]				
Finland	0.054 [0.007]	0.052 [0.005]	0.048 [0.008]	0.020 [0.006]	0.028 [0.003]				
France	0.084 [0.003]	0.086 [0.003]	0.072 [0.005]	0.090 [0.004]	0.098 [0.003]				
Germany	0.069 [0.004]	0.059 [0.003]	0.047 [0.004]	0.070 [0.004]	0.079 [0.003]				
Greece	0.144 [0.010]	0.121 [0.006]	0.117 [0.009]	0.126 [0.008]	0.109 [0.008]				
Hungary	0.104 [0.003]	0.103 [0.003]	0.095 [0.004]	0.113 [0.003]	0.108 [0.003]				
Iceland	0.061 [0.018]	0.032 [0.006]	0.030 [0.006]	0.012 [0.004]	0.016 [0.003]				
Ireland	0.102 [0.006]	0.097 [0.005]	0.048 [0.010]	0.084 [0.007]	0.078 [0.005]				
Italy	0.116 [0.004]	0.091 [0.002]	0.080 [0.006]	0.108 [0.003]	0.097 [0.004]				
Latvia	0.127 [0.008]	0.110 [0.005]	0.095 [0.008]	0.110 [0.007]	0.111 [0.006]				
Lithuania	0.088 [0.008]	0.067 [0.005]	0.059 [0.006]	0.069 [0.008]	0.067 [0.006]				
Luxembourg	0.132 [0.005]	0.121 [0.003]	0.090 [0.009]	0.133 [0.004]	0.136 [0.003]				
Malta	0.085 [0.006]	0.080 [0.005]	0.057 [0.005]	0.071 [0.006]	0.072 [0.005]				
Netherlands	0.064 [0.007]	0.053 [0.003]	0.041 [0.006]	0.028 [0.004]	0.019 [0.002]				
Norway	0.041 [0.007]	0.041 [0.005]	0.030 [0.006]	0.020 [0.004]	0.023 [0.003]				
Poland	0.104 [0.003]	0.097 [0.003]	0.095 [0.004]	0.102 [0.004]	0.099 [0.004]				
Portugal	0.133 [0.005]	0.124 [0.005]	0.116 [0.007]	0.136 [0.006]	0.127 [0.007]				
Romania	0.161 [0.006]	0.104 [0.005]	0.119 [0.007]	0.120 [0.006]	0.111 [0.006]				
Slovakia	0.058 [0.004]	0.051 [0.003]	0.042 [0.005]	0.050 [0.004]	0.046 [0.004]				
Slovenia	0.074 [0.004]	0.073 [0.004]	0.059 [0.004]	0.032 [0.004]	0.036 [0.002]				
Spain	0.141 [0.005]	0.120 [0.003]	0.089 [0.010]	0.128 [0.003]	0.120 [0.008]				
Sweden	0.089 [0.036]	0.060 [0.009]	0.025 [0.009]	0.025 [0.004]	0.031 [0.003]				
Switzerland	0.091 [0.005]	0.083 [0.004]	0.063 [0.008]	0.080 [0.006]	0.090 [0.004]				
United Kingdom	0.093 [0.008]	0.090 [0.005]	0.062 [0.010]	0.071 [0.008]	0.079 [0.004]				

 TABLE D.2 – Inequality of Opportunity Estimates

Data: EU-SILC 2011 cross-sectional (rev.5, June 2015). **Note:** This tables shows inequality of opportunity estimates for different estimation methods. Inequality of opportunity is measured by the Gini coefficient of the counterfactual distribution \hat{y}^{C} . Standard errors are derived from 200 bootstrapped re-samples of the test data.

The Roots of Inequality: Estimating Inequality of Opportunity from Regression Trees and Forests

Supplementary Material

S.1 DESCRIPTIVE STATISTICS

		Birth Area		Parents	s in HH	HF			
Country	Male	Native	EU	Both	None	Adults	Working Ad.	Children	Home Owner
Austria	0.501	0.790	0.070	0.856	0.017	2.730	1.760	2.600	0.585
Belgium	0.498	0.824	0.076	0.855	0.019	2.380	1.590	2.780	0.750
Bulgaria	0.500	0.994	0.001	0.904	0.012	2.440	2.010	2.070	0.910
Croatia	0.501	0.875	0.017	0.874	0.020	2.560	1.350	2.310	0.902
Cyprus	0.525	0.787	0.096	0.900	0.015	2.640	1.670	2.700	0.784
Czech Republic	0.508	0.964	0.026	0.851	0.013	2.090	1.920	2.240	0.597
Denmark	0.505	0.923	0.026	0.809	0.027	2.220	2.310	2.240	0.736
Estonia	0.525	0.847	0.000	0.756	0.011	2.100	1.800	2.090	0.859
Finland	0.499	0.954	0.018	0.829	0.016	2.360	1.750	2.300	0.772
France	0.509	0.885	0.036	0.820	0.022	2.470	1.660	1.750	0.630
Germany	0.496	0.868	0.000	0.830	0.020	2.240	1.680	2.320	0.499
Greece	0.498	0.890	0.025	0.931	0.019	2.310	1.560	2.330	0.834
Hungary	0.517	0.988	0.008	0.844	0.041	2.140	1.750	2.270	0.830
Iceland	0.507	0.920	0.042	0.899	0.012	2.420	1.900	2.630	0.893
Ireland	0.524	0.783	0.149	0.893	0.078	3.170	3.200	3.200	0.727
Italy	0.502	0.880	0.040	0.901	0.011	2.590	1.620	2.410	0.685
Latvia	0.520	0.865	0.000	0.763	0.012	1.970	1.760	2.280	0.455
Lithuania	0.521	0.939	0.004	0.846	0.016	2.320	2.020	2.460	0.698
Luxembourg	0.499	0.480	0.401	0.868	0.020	2.530	1.640	2.710	0.734
Malta	0.497	0.944	0.000	0.932	0.020	3.020	1.840	2.680	0.576
Netherlands	0.509	0.903	0.020	0.882	0.016	2.100	1.540	3.250	0.575
Norway	0.511	0.907	0.041	0.913	0.014	2.020	1.760	1.870	0.922
Poland	0.504	0.999	0.000	0.889	0.015	2.700	1.960	2.440	0.644
Portugal	0.506	0.906	0.022	0.854	0.017	2.680	2.230	2.680	0.544
Romania	0.506	0.999	0.000	0.919	0.009	2.770	1.900	2.270	0.861
Slovakia	0.519	0.987	0.010	0.920	0.010	2.520	2.080	2.340	0.694
Slovenia	0.496	0.876	0.000	0.855	0.019	2.530	1.770	2.200	0.746
Spain	0.495	0.834	0.051	0.893	0.012	2.880	2.110	2.430	0.819
Sweden	0.493	0.846	0.050	0.820	0.035	2.070	1.780	2.350	0.757
Switzerland	0.505	0.684	0.197	0.837	0.017	2.550	1.900	2.530	0.546
United Kingdom	0.507	0.848	0.042	0.825	0.024	2.340	2.240	2.410	0.649

TABLE S.1 – Descriptive Statistics (Individual and Household)

Data: EU-SILC 2011 cross-sectional (rev.5, June 2015). **Note:** Omitted circumstance expressions listed in order of the circumstance categories: "Female"; "Non-EU"; "Only Mother/Only Father/Collective House"; "House Not Owned". See also Table 1.

	Birth Area		Birth Area		Birth Area		Citize	enship		Educatio	1		Act	ivity					Occupatio	on (ISCO-	08 1-Digit	t)			
Country	Native	EU	Resid.	EU	Low	Med.	High	Empl.	Self- Empl.	Un- empl.	House Work	1	2	3	4	5	6	7	8	9	Superv.				
Austria	0.743	0.093	0.777	0.068	0.398	0.421	0.135	0.714	0.215	0.003	0.001	0.043	0.046	0.064	0.051	0.138	0.145	0.284	0.063	0.085	0.338				
Belgium	0.748	0.100	0.762	0.093	0.491	0.199	0.178	0.699	0.179	0.007	0.002	0.068	0.126	0.104	0.084	0.054	0.057	0.209	0.127	0.041	0.278				
Bulgaria	0.933	0.004	0.936	0.001	0.466	0.333	0.081	0.899	0.028	0.005	0.000	0.022	0.065	0.047	0.029	0.035	0.135	0.216	0.207	0.142	0.093				
Croatia	0.822	0.006	0.834	0.004	0.464	0.312	0.063	0.763	0.103	0.037	0.016	0.025	0.041	0.088	0.036	0.072	0.049	0.214	0.103	0.228	0.129				
Cyprus	0.803	0.082	0.808	0.094	0.667	0.178	0.091	0.566	0.381	0.004	0.000	0.011	0.071	0.074	0.029	0.104	0.161	0.245	0.122	0.125	0.229				
Czech Republic	0.878	0.065	0.910	0.036	0.602	0.216	0.090	0.891	0.017	0.001	0.000	0.033	0.070	0.125	0.036	0.035	0.039	0.305	0.195	0.053	0.233				
Denmark	0.935	0.025	0.970	0.020	0.368	0.418	0.214	0.708	0.272	0.004	0.001	0.111	0.122	0.070	0.043	0.103	0.160	0.288	0.072	0.009	0.447				
Estonia	0.603	0.270	0.637	0.233	0.300	0.338	0.165	0.823	0.006	0.003	0.002	0.076	0.092	0.053	0.014	0.013	0.034	0.221	0.253	0.053	0.153				
Finland	0.827	0.007	0.827	0.007	0.491	0.182	0.162	0.592	0.209	0.016	0.001	0.041	0.089	0.085	0.016	0.046	0.135	0.146	0.138	0.044					
France	0.789	0.078	0.857	0.057	0.695	0.073	0.095	0.753	0.170	0.003	0.001	0.084	0.068	0.111	0.072	0.038	0.103	0.155	0.055	0.223	0.335				
Germany	0.800	0.200	0.855	0.145	0.125	0.496	0.213	0.819	0.123	0.008	0.001	0.046	0.104	0.158	0.051	0.061	0.059	0.266	0.154	0.040	0.299				
Greece	0.887	0.016	0.911	0.015	0.587	0.135	0.084	0.449	0.517	0.002	0.000	0.073	0.047	0.026	0.087	0.046	0.308	0.210	0.099	0.055	0.182				
Hungary	0.962	0.017	0.969	0.012	0.599	0.241	0.087	0.892	0.043	0.001	0.001	0.037	0.060	0.052	0.017	0.053	0.094	0.279	0.193	0.137	0.117				
Iceland	0.918	0.050	0.923	0.044	0.334	0.486	0.139	0.638	0.332	0.001	0.000	0.115	0.121	0.076	0.024	0.094	0.180	0.220	0.094	0.042	0.570				
Ireland	0.792	0.107	0.758	0.094	0.574	0.258	0.112	0.659	0.221	0.049	0.002	0.104	0.092	0.042	0.022	0.072	0.155	0.149	0.065	0.158	0.344				
Italy	0.823	0.022	0.827	0.020	0.708	0.136	0.038	0.614	0.244	0.016	0.004	0.054	0.040	0.074	0.057	0.068	0.099	0.227	0.105	0.118	0.199				
Latvia	0.572	0.248	0.642	0.165	0.381	0.297	0.098	0.767	0.005	0.002	0.003	0.036	0.083	0.037	0.010	0.019	0.069	0.199	0.218	0.083	0.070				
Lithuania	0.899	0.004	0.926	0.004	0.538	0.228	0.085	0.916	0.011	0.000	0.001	0.049	0.074	0.038	0.017	0.023	0.080	0.241	0.179	0.214	0.110				
Luxembourg	0.387	0.467	0.400	0.466	0.484	0.316	0.120	0.757	0.174	0.001	0.001	0.063	0.093	0.118	0.048	0.035	0.112	0.228	0.183	0.039	0.251				
Malta	0.952	0.041	0.953	0.040	0.561	0.180	0.059	0.717	0.214	0.013	0.001	0.062	0.046	0.106	0.045	0.141	0.050	0.244	0.099	0.106	0.225				
Netherlands	0.829	0.028	0.888	0.022	0.376	0.285	0.198	0.726	0.173	0.006	0.006	0.087	0.124	0.155	0.051	0.069	0.086	0.200	0.079	0.031	0.310				
Norway	0.897	0.046	0.908	0.041	0.328	0.390	0.278	0.712	0.255	0.002	0.001	0.116	0.110	0.167	0.029	0.057	0.111	0.227	0.100	0.032	0.285				
Poland	0.955	0.012	0.980	0.003	0.462	0.448	0.070	0.701	0.238	0.002	0.001	0.036	0.044	0.053	0.025	0.042	0.237	0.254	0.157	0.078	0.111				
Portugal	0.932	0.006	0.945	0.006	0.700	0.031	0.031	0.650	0.248	0.002	0.001	0.047	0.032	0.060	0.038	0.082	0.185	0.264	0.114	0.077	0.190				
Romania	0.938	0.001	0.939	0.001	0.726	0.088	0.030	0.642	0.237	0.004	0.013	0.004	0.040	0.034	0.016	0.018	0.253	0.249	0.121	0.104	0.045				
Slovakia	0.935	0.020	0.945	0.011	0.362	0.497	0.075	0.921	0.011	0.002	0.001	0.042	0.060	0.095	0.028	0.043	0.030	0.285	0.209	0.128	0.145				
Slovenia	0.769	0.200	0.000	0.000	0.684	0.166	0.085	0.773	0.099	0.013	0.011	0.024	0.052	0.100	0.037	0.052	0.089	0.257	0.080	0.173	0.242				
Spain	0.836	0.047	0.846	0.046	0.762	0.064	0.081	0.702	0.219	0.006	0.001	0.056	0.045	0.076	0.055	0.087	0.145	0.191	0.113	0.137	0.191				
Sweden	0.945	0.022	0.851	0.061	0.422	0.350	0.182	0.745	0.211	0.002	0.001	0.043	0.118	0.067	0.031	0.092	0.086	0.230	0.108	0.019	0.337				
Switzerland	0.588	0.286	0.603	0.280	0.227	0.487	0.151	0.653	0.292	0.001	0.000	0.086	0.131	0.140	0.057	0.060	0.111	0.223	0.077	0.054	0.397				
United Kingdom	0.800	0.064	0.869	0.039	0.508	0.228	0.150	0.795	0.147	0.025	0.002	0.095	0.142	0.085	0.040	0.075	0.036	0.236	0.133	0.083	0.398				

TABLE S.2 – Descriptive Statistics (Fathers)

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

Note: Omitted circumstance expressions listed in order of the circumstance categories: "Non-EU"; "Not Europe"; "Dead/Unknown/Illiterate"; "Dead/Unknown/Retired/Other Inactive"; "Dead/Unknown/Not Working/Armed Forces"; "Dead/Unknown/Not Working/Non-Supervisory". See also Table 1.

	Birth Area		Citizenship		Education			Activity				Occupation (ISCO-08 1-Digit)									
Country	Native	EU	Resid.	EU	Low	Med.	High	Empl.	Self- Empl.	Un- empl.	House Work	1	2	3	4	5	6	7	8	9	Superv.
Austria	0.740	0.096	0.789	0.065	0.587	0.328	0.041	0.714	0.215	0.003	0.001	0.009	0.024	0.009	0.071	0.154	0.128	0.045	0.010	0.087	0.092
Belgium	0.755	0.097	0.790	0.092	0.564	0.201	0.126	0.699	0.179	0.007	0.002	0.009	0.081	0.045	0.058	0.046	0.002	0.016	0.024	0.069	0.034
Bulgaria	0.931	0.003	0.981	0.002	0.464	0.357	0.099	0.899	0.028	0.005	0.000	0.008	0.123	0.040	0.092	0.140	0.181	0.099	0.064	0.152	0.030
Croatia	0.823	0.008	0.848	0.003	0.634	0.189	0.043	0.763	0.103	0.037	0.016	0.003	0.058	0.036	0.046	0.070	0.022	0.034	0.013	0.122	0.033
Cyprus	0.804	0.080	0.812	0.091	0.684	0.162	0.057	0.566	0.381	0.004	0.000	0.001	0.045	0.020	0.037	0.067	0.036	0.022	0.042	0.220	0.048
Czech Republic	0.882	0.061	0.946	0.037	0.670	0.261	0.043	0.891	0.017	0.001	0.000	0.011	0.080	0.104	0.149	0.160	0.074	0.105	0.080	0.139	0.088
Denmark	0.922	0.029	0.935	0.023	0.531	0.283	0.186	0.708	0.272	0.004	0.001	0.020	0.103	0.095	0.123	0.225	0.035	0.052	0.026	0.001	0.122
Estonia	0.601	0.272	0.726	0.250	0.334	0.391	0.208	0.823	0.006	0.003	0.002	0.049	0.169	0.109	0.097	0.110	0.084	0.051	0.124	0.113	0.085
Finland	0.826	0.007	0.933	0.006	0.559	0.238	0.145	0.592	0.209	0.016	0.001	0.012	0.126	0.091	0.122	0.145	0.048	0.046	0.057	0.202	
France	0.806	0.067	0.880	0.047	0.724	0.079	0.072	0.753	0.170	0.003	0.001	0.011	0.036	0.050	0.111	0.107	0.059	0.049	0.005	0.109	0.072
Germany	0.811	0.189	0.862	0.138	0.284	0.475	0.081	0.819	0.123	0.008	0.001	0.012	0.051	0.079	0.089	0.116	0.025	0.015	0.087	0.033	0.059
Greece	0.888	0.016	0.916	0.016	0.592	0.133	0.044	0.449	0.517	0.002	0.000	0.023	0.027	0.004	0.039	0.048	0.223	0.034	0.021	0.049	0.026
Hungary	0.964	0.016	0.980	0.012	0.655	0.243	0.053	0.892	0.043	0.001	0.001	0.014	0.049	0.063	0.113	0.118	0.061	0.075	0.087	0.167	0.044
Iceland	0.905	0.059	0.924	0.046	0.626	0.275	0.075	0.638	0.332	0.001	0.000	0.030	0.095	0.045	0.109	0.180	0.064	0.028	0.013	0.130	0.149
Ireland	0.787	0.114	0.761	0.103	0.546	0.324	0.097	0.659	0.221	0.049	0.002	0.022	0.061	0.007	0.052	0.059	0.017	0.014	0.007	0.060	0.082
Italy	0.820	0.024	0.862	0.024	0.779	0.112	0.023	0.614	0.244	0.016	0.004	0.011	0.038	0.022	0.029	0.051	0.035	0.031	0.022	0.062	0.041
Latvia	0.585	0.234	0.793	0.182	0.414	0.399	0.125	0.767	0.005	0.002	0.003	0.031	0.138	0.084	0.098	0.121	0.085	0.093	0.023	0.221	0.074
Lithuania	0.902	0.002	0.959	0.003	0.519	0.316	0.106	0.916	0.011	0.000	0.001	0.035	0.129	0.046	0.049	0.109	0.067	0.112	0.034	0.293	0.068
Luxembourg	0.374	0.483	0.393	0.485	0.587	0.245	0.071	0.757	0.174	0.001	0.001	0.028	0.049	0.046	0.036	0.061	0.054	0.015	0.024	0.108	0.047
Malta	0.950	0.043	0.957	0.038	0.652	0.145	0.026	0.717	0.214	0.013	0.001	0.003	0.019	0.007	0.009	0.018	0.002	0.004	0.009	0.010	0.011
Netherlands	0.829	0.027	0.907	0.023	0.532	0.288	0.087	0.726	0.173	0.006	0.006	0.010	0.050	0.038	0.052	0.089	0.016	0.011	0.008	0.060	0.037
Norway	0.877	0.048	0.891	0.043	0.368	0.437	0.181	0.712	0.255	0.002	0.001	0.031	0.041	0.142	0.114	0.209	0.053	0.017	0.026	0.091	0.065
Poland	0.957	0.010	0.990	0.004	0.524	0.410	0.057	0.701	0.238	0.002	0.001	0.018	0.057	0.053	0.071	0.096	0.262	0.080	0.018	0.118	0.050
Portugal	0.928	0.008	0.950	0.007	0.631	0.029	0.028	0.650	0.248	0.002	0.001	0.016	0.031	0.017	0.025	0.075	0.158	0.059	0.032	0.145	0.048
Romania	0.936	0.001	0.939	0.001	0.728	0.112	0.020	0.642	0.237	0.004	0.013	0.001	0.034	0.024	0.026	0.050	0.218	0.076	0.040	0.080	0.010
Slovakia	0.932	0.023	0.980	0.010	0.451	0.482	0.039	0.921	0.011	0.002	0.001	0.010	0.075	0.110	0.107	0.161	0.034	0.096	0.052	0.203	0.048
Slovenia	0.791	0.178	0.000	0.000	0.752	0.148	0.058	0.773	0.099	0.013	0.011	0.006	0.047	0.093	0.085	0.090	0.061	0.066	0.006	0.193	0.089
Spain	0.836	0.046	0.849	0.046	0.802	0.048	0.040	0.702	0.219	0.006	0.001	0.010	0.025	0.010	0.021	0.059	0.028	0.021	0.009	0.071	0.029
Sweden	0.942	0.024	0.855	0.058	0.409	0.369	0.201	0.745	0.211	0.002	0.001	0.006	0.087	0.033	0.057	0.152	0.016	0.009	0.021	0.035	0.095
Switzerland	0.567	0.307	0.599	0.286	0.410	0.399	0.057	0.653	0.292	0.001	0.000	0.027	0.056	0.069	0.069	0.125	0.055	0.039	0.025	0.068	0.064
United Kingdom	0.808	0.064	0.877	0.036	0.679	0.099	0.124	0.795	0.147	0.025	0.002	0.026	0.097	0.068	0.078	0.152	0.005	0.028	0.044	0.127	0.104

 TABLE S.3 – Descriptive Statistics (Mothers)

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

Note: Omitted circumstance expressions listed in order of the circumstance categories are: "Non-EU"; "Not Europe"; "Dead/Unknown/Illiterate"; "Dead/Unknown/Retired/Other Inactive"; "Dead/Unknown/Not Working/Armed Forces"; "Dead/Unknown/Not Working/Non-Supervisory". See also Table 1.

S.2 ALTERNATIVE INEQUALITY INDEXES



FIGURE S.1 – Correlation of Estimates by Method (GE[0])

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

Note: This figure shows inequality of opportunity estimates from different estimation methods relative to random forests. Inequality of opportunity is measured by the general entropy measure ($\alpha = 0$) of the counterfactual distribution \hat{y}^{C} .



FIGURE S.2 – Correlation of Estimates by Method (GE[1])

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

Note: This figure shows inequality of opportunity estimates from different estimation methods relative to random forests. Inequality of opportunity is measured by the general entropy measure ($\alpha = 1$) of the counterfactual distribution \hat{y}^{C} .



FIGURE S.3 – Correlation of Estimates by Method (GE[2])

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

Note: This figure shows inequality of opportunity estimates from different estimation methods relative to random forests. Inequality of opportunity is measured by the general entropy measure ($\alpha = 2$) of the counterfactual distribution \hat{y}^{C} .

S.3 OPPORTUNITY STRUCTURES



FIGURE S.4 – Opportunity Tree (Austria)

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



FIGURE S.5 – Opportunity Tree (Belgium)

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



FIGURE S.6 – Opportunity Tree (Croatia)

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



FIGURE S.7 – Opportunity Tree (Cyprus)





FIGURE S.8 – Opportunity Tree (Czech Republic)



FIGURE S.9 – Opportunity Tree (Denmark)



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



FIGURE S.10 – Opportunity Tree (Estonia)

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

FIGURE S.11 – Opportunity Tree (Finland)



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



FIGURE S.12 – Opportunity Tree (France)

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

FIGURE S.13 – Opportunity Tree (Germany)



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

FIGURE S.14 – Opportunity Tree (Greece)



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

FIGURE S.15 – Opportunity Tree (Hungary)



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

FIGURE S.16 – Opportunity Tree (Iceland)



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



FIGURE S.17 – Opportunity Tree (Ireland)

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



FIGURE S.18 – Opportunity Tree (Italy)

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

FIGURE S.19 – Opportunity Tree (Latvia)



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



FIGURE S.20 – Opportunity Tree (Lithuania)

Data: EU-SILC 2011 cross-sectional (rev.5, June 2015).
FIGURE S.21 – Opportunity Tree (Luxembourg)



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

FIGURE S.22 – Opportunity Tree (Malta)



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



FIGURE S.23 – Opportunity Tree (Netherlands)

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

FIGURE S.24 – Opportunity Tree (Norway)



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



FIGURE S.25 – Opportunity Tree (Poland)

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

FIGURE S.26 – Opportunity Tree (Portugal)



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

FIGURE S.27 – Opportunity Tree (Romania)



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

FIGURE S.28 – Opportunity Tree (Slovakia)



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



FIGURE S.29 – Opportunity Tree (Slovenia)

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

FIGURE S.30 – Opportunity Tree (Spain)



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

FIGURE S.31 – Opportunity Tree (Sweden)



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



FIGURE S.32 – Opportunity Tree (Switzerland)

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

FIGURE S.33 – Opportunity Tree (United Kingdom)



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

Note: The tree is constructed by the conditional inference algorithm (Section 3.1). The set of observed circumstances Ω used to construct the conditional inference tree is detailed in Table 1. Ellipses indicate splitting points, while the rectangular boxes indicate terminal nodes. Within each ellipse we indicate the splitting variable as well as the *p*-value associated with the respective split. The first number inside the terminal nodes indicates the population share belonging to the circumstance type, while the second number shows the respective estimate of the conditional expectation y^{C} .



FIGURE S.34 – Variable Importance Plot from Forests

Note: Each dot shows the importance of a particular circumstance variable ω^p . Variable importance is measured by the decrease in MSE^{OOB} after permuting ω^p such that it is orthogonal to *y*. The importance measure is standardized such that the circumstance with the greatest importance in each country equals 1.

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