

# The Roots of Inequality

## Estimating Inequality of Opportunity from Regression Trees and Forests

Conference on Inequality of Opportunity  
University of Queensland & Griffith University  
27-28 June, 2019

Paolo Brunori  
*University of Florence & University of Bari*

- “The Roots of Inequality Estimating Inequality of Opportunity from Regression Trees and Forests” is a joint work with Paul Hufe (Ifo) and Daniel G. Mahler (World Bank);

## Ex-ante IOP estimation

- $y = g(C) + \epsilon$
- causality is conceptually and empirically excluded  $\rightarrow$  covariance of the outcome and circumstances' variability

$$IOP = I(\hat{y})$$

$I$  is a suitable inequality index;

$\hat{y}$  is the predicted outcome distribution based on  $\hat{g}(C)$ ;

$\hat{g}()$  is estimated on survey data.

# Typical machine learning domain

- unknown data generating process;
- need to establish a reliable empirical link between a set of controls and an outcome.

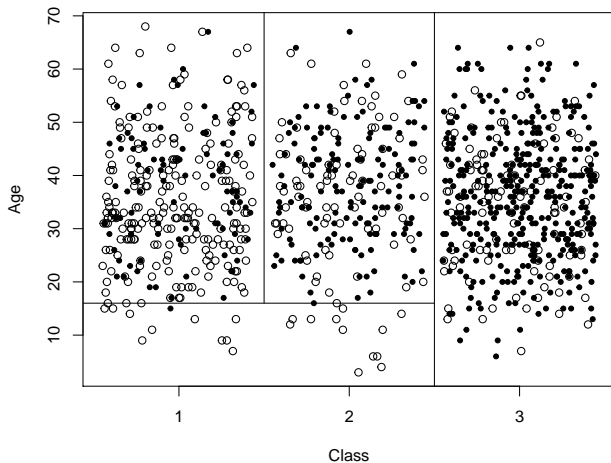
# ML and IOP

- ML: bias-variance trade-off;
- IOP partial observability (downward bias) - sampling variance (upward bias);
- ML: choose the model that minimizes out-of-sample MSE;
- IOP: choose the model that maximizes IOP out-of-sample (Social Choice and Welfare - 2019 with Peragine and Serlenga).

# Trees

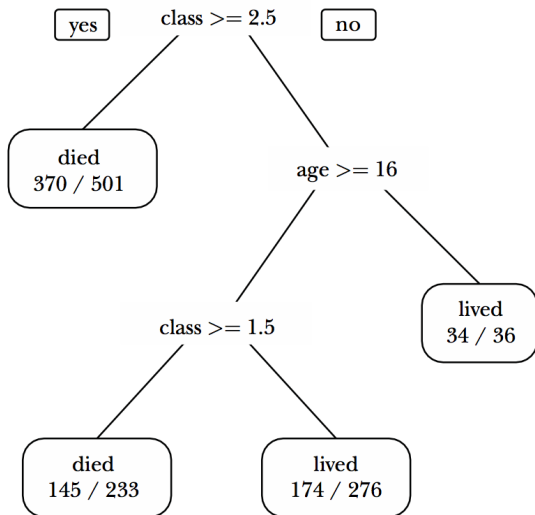
- among supervised learning algorithms regression trees seem an obvious choice;
- a tree is an algorithm to predict a dependent variable based on observable predictors (Morgan and Sonquist,1963; Breiman et al.,1984);
- the population is divided into non-overlapping subgroups based on a partition of the predictors' space;
- prediction of each observation is the the mean value of the dependent variable in the group.

# What is a tree? cnt.



*Source: adapted from Varian, 2014*

What is a tree? cnt.



Source: Varian, 2014



# Tuning

- a very deep tree performs poorly out-of-sample;
- different solutions to prevent overfitting lead to different type of trees;
- *conditional inference trees* condition each split on a statistical test (Hothorn et al., 2006).

## Conditional inference trees

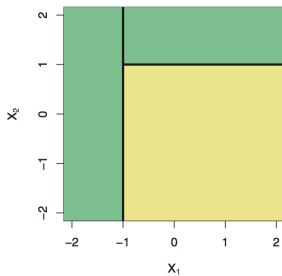
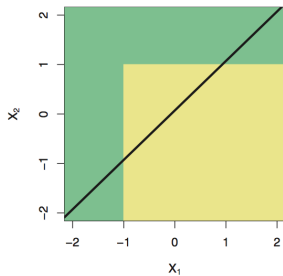
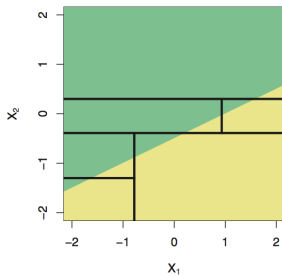
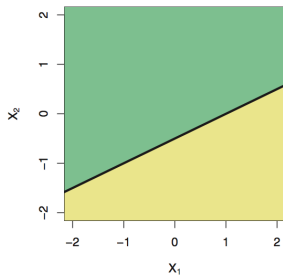
- test the null hypothesis of independence,  
 $H^{C_p} = D(Y|C_p) = D(Y), \forall C_p \in \mathbf{C}$ ;
- no (adjusted) p-value  $< \alpha \rightarrow$  exit the algorithm;
- select the variable,  $C^*$ , with the lowest p-value;
- test the discrepancy between the subsamples for each possible binary partition based on  $C^*$ ;
- split the sample by selecting the splitting point that yields the lowest p-value;
- repeat the algorithm for each of the resulting subsamples.

## Opportunity trees: *pros*

- the selection of  $\mathbf{C}$  is no longer arbitrary;
- the model specification is endogenous to data;
- provide a test for the null hypothesis of *EOP*;
- tell a story about the opportunity structure.

## Opportunity trees: *cons*

- misleading when two or more controls are highly correlated;
- perform poorly if the data generating process is linear.



source: James et al. (2013)

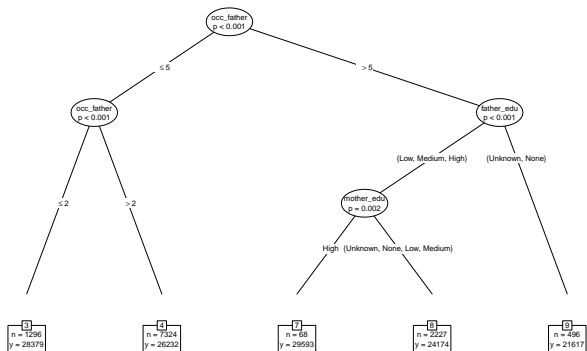
# Random forests

- random forests improve tree's predictive performance;
- a forest is made of hundreds of conditional inference trees;
- each tree uses a subsample of observations and, at each splitting point, a subsample of controls.

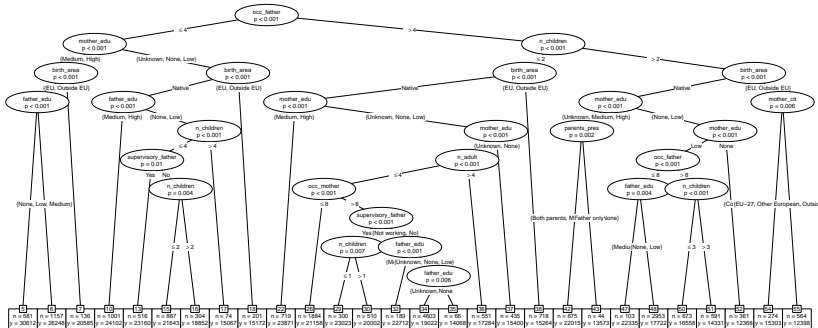
## data

- EU-SILC 2011;
- subsample: adults (30-60);
- $y$ : household equivalent disposable income;
- **C**: 21 questions about respondents' background (sex, birth area, proxies for socioeconomic status);
- already used to estimate IOP.

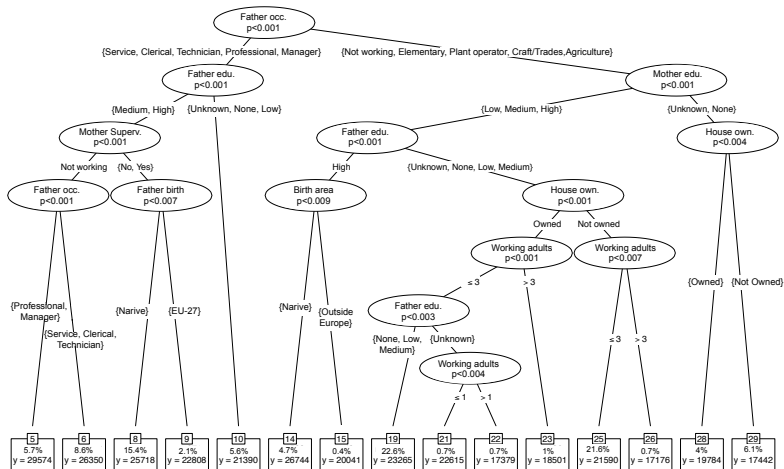
# The Netherlands







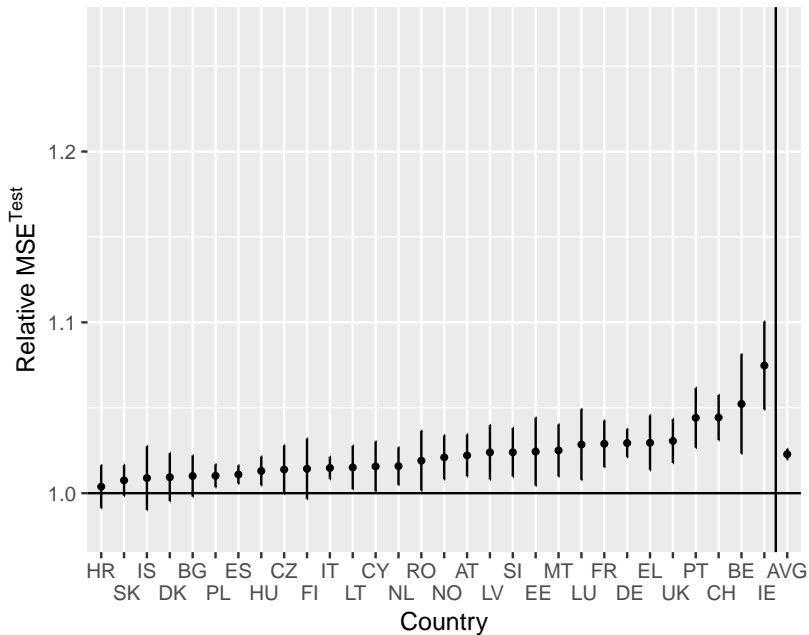
# Germany



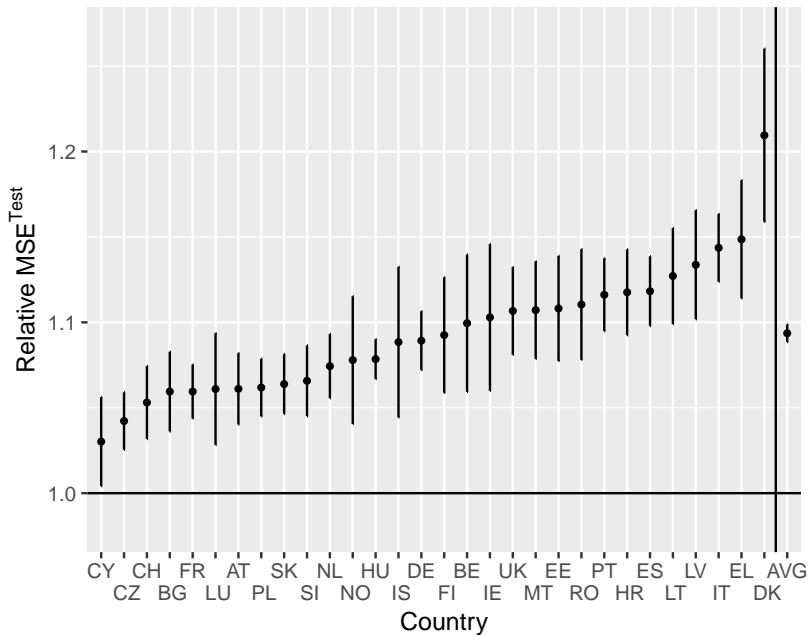
# Random forests

- random forests of 200 conditional inference trees used to:
  - estimate IOp;
  - quantify relative variable importance.

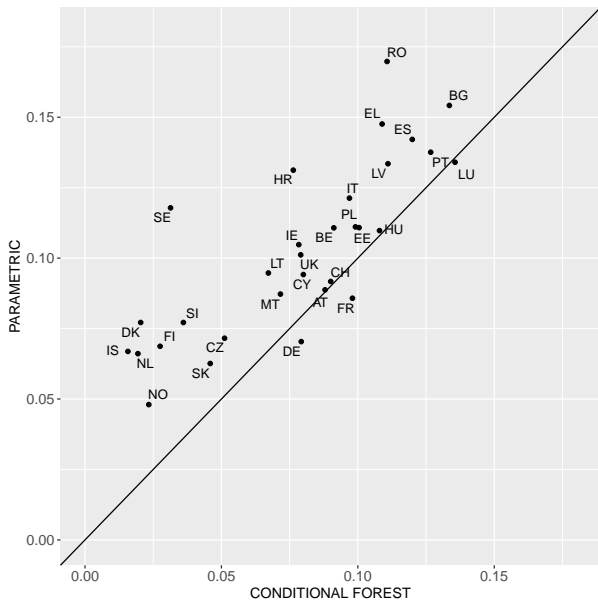
# Predictive performance: trees Vs. forest



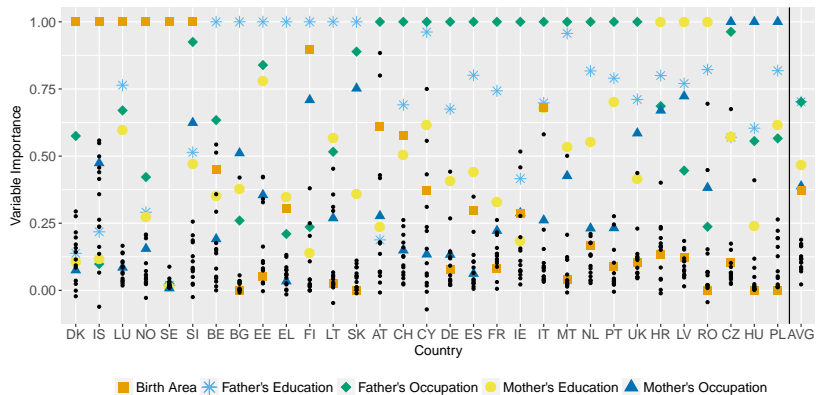
# Predictive performance: Parametric Vs. forest



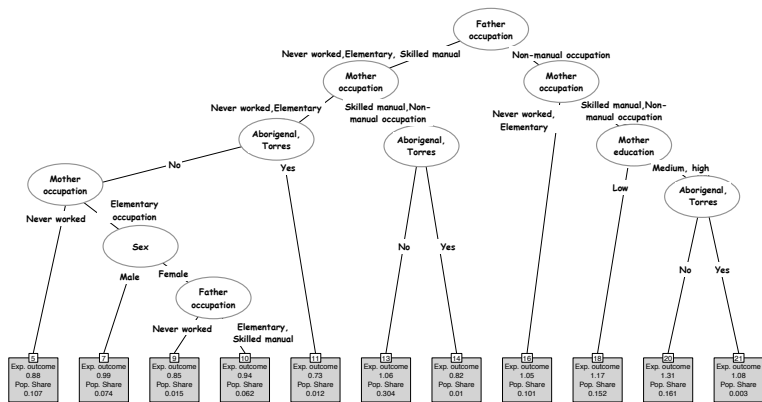
# Estimates: Parametric Vs. forest



# Variables importance



# Bonus tree: Australia, 2015





# Ex-post IOP

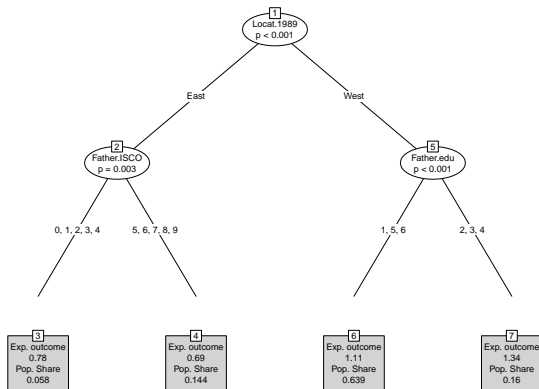
(joint work in progress with Guido Neidhöfer)

- Conditional inference regression trees have two important advantage
  - they identify types;
  - they are parsimonious.
- having types with sufficient sample size one can move further and estimate IOP consistently with Roemer's original theory;
- ex-post IOP definition is based on the estimation of the type-specific outcome distribution.

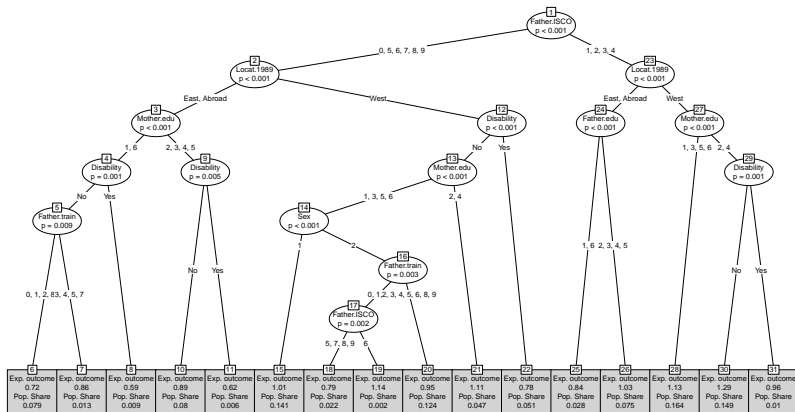
# Effort

- According to Roemer the quantile of the type-specific outcome distribution is a convincing proxy of the degree of effort exerted;
- ex-post IOP quantifies to what extent individuals exerting the same degree of effort do not obtain the same outcome;
- we use Bernstein polynomial approximation of the types' ECDF to measure ex-post IOP.

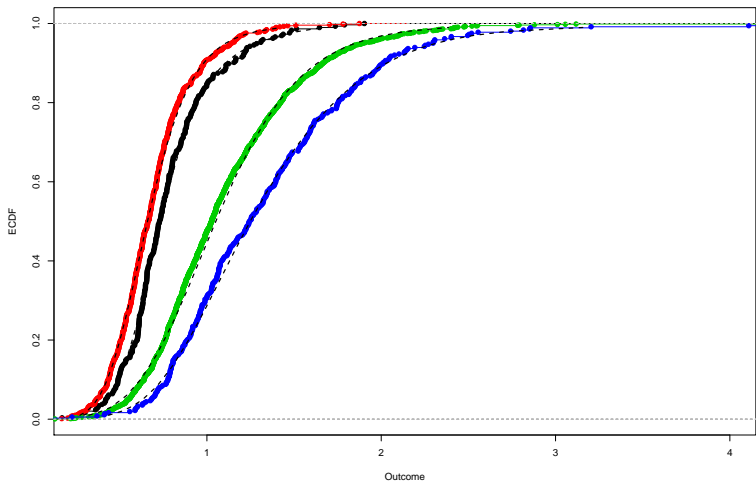
# Opportunity tree in 1992



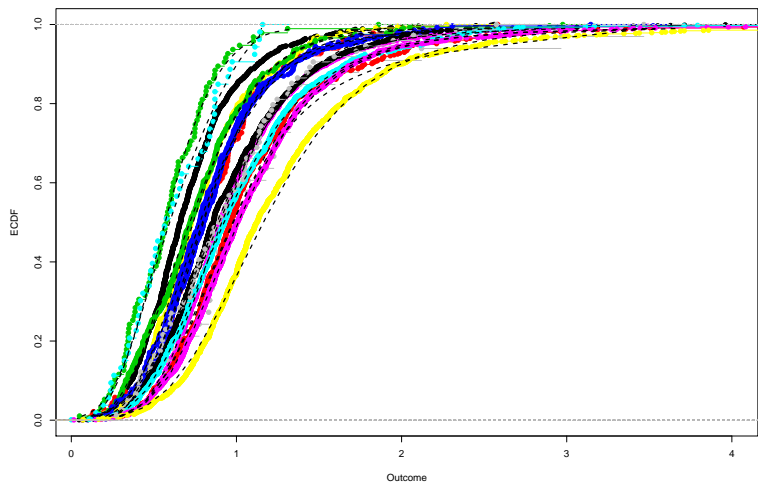
# Opportunity tree in 2016



# IOP in 1992



# IOP in 2016



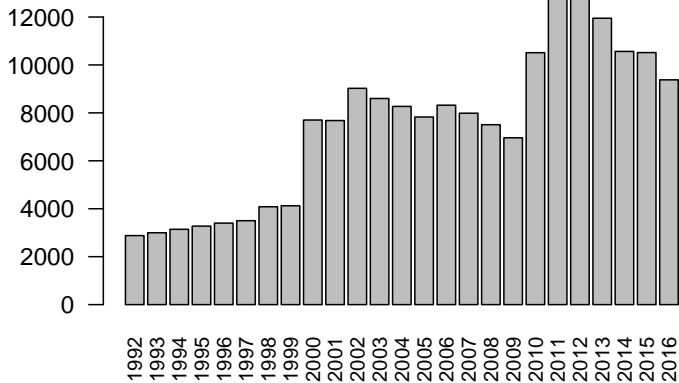
# Conclusions

- many other ML approaches can be used:
  - unsupervised learning such as Li Donni et al. (2015) and Wu, Trivedi, Rao, Tang (2018)
  - best subset regression (EqualChances.org)
  - LASSO (or other regularization methods) as for example Hufe et al. (2019)
- but there exists a second key trade off in ML: complexity Vs. interpretability.

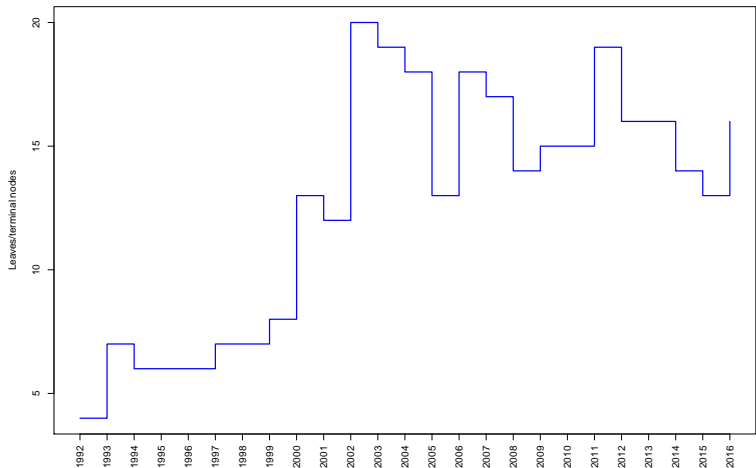
Additional material: trend in Germany



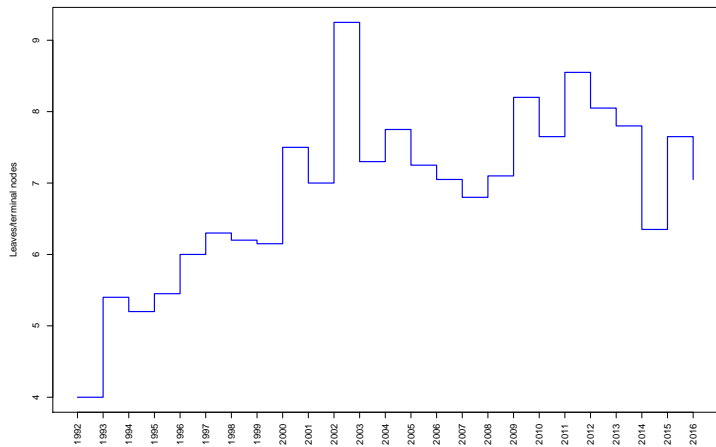
# Sample size 1992-2016



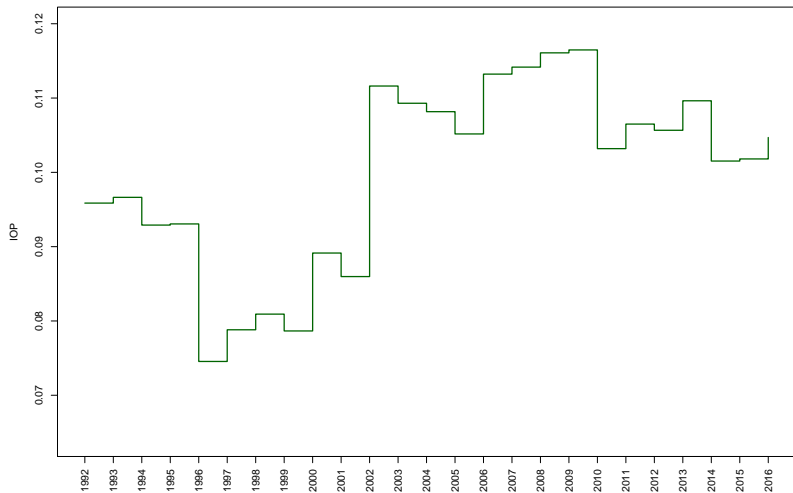
# Number of types 1992-2016



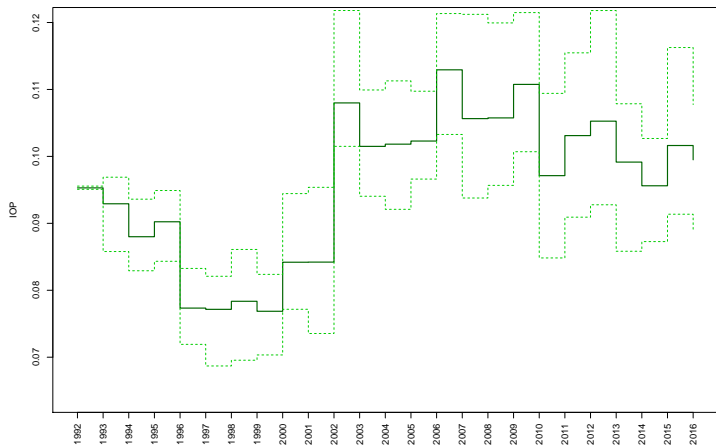
# Mean number of types (same sample size) 1992-2016



# IOP trend 1992-2016



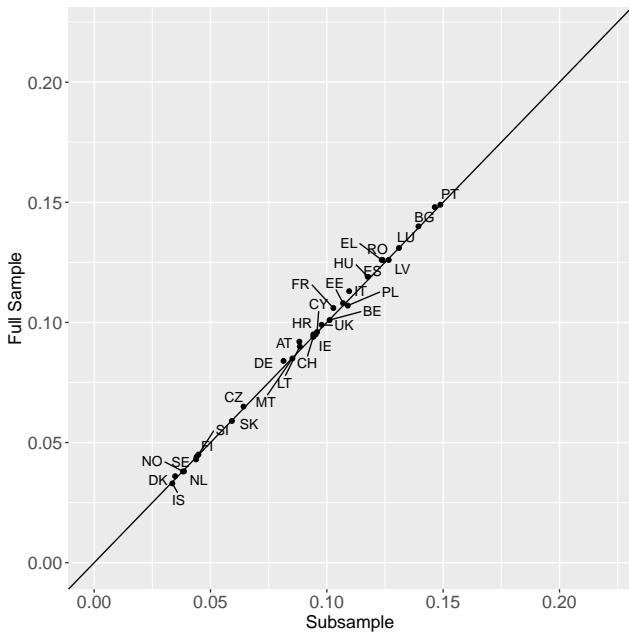
## Mean IOP trend 1992-2016 (same sample size)



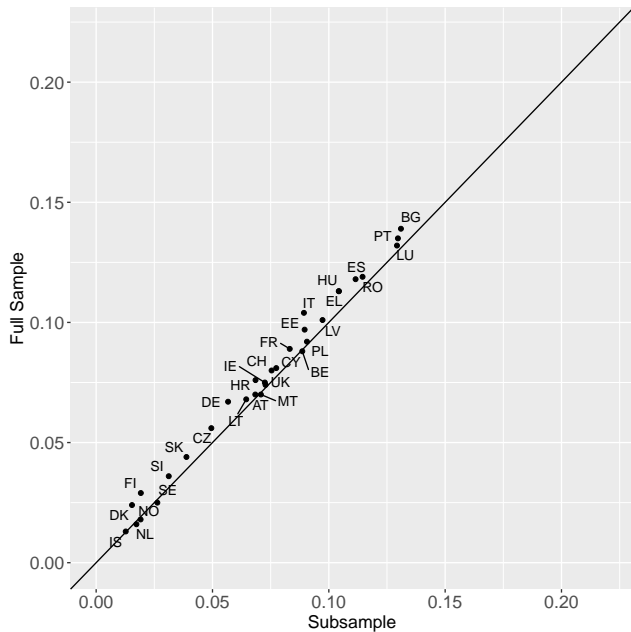
*Confidence bounds are the 0.975 and 0.025 quantiles of the distribution of IOP estimates.*

# Additional material: sample size EU-SILC

## Sensitivity to sample size: forests



## Sensitivity to sample size: trees





## Sensitivity to sample size: parametric

