

Modified Causal Forest: Estimation

Section 1: General information

Welcome to the mcf estimation and optimal policy package.

This report provides you with a summary of specifications and results. More detailed information can be found in the respective output files. Figures and data (in csv-format, partly to recreate the figures on your own) are provided in the output path as well.

Output information for MCF ESTIMATION

All outputs: /Users/zg/CML_toy/output

Subdirectories with figures and data are named ate_iate, gate, and common support and contain the content related to their name.

Detailed text output: /Users/zg/CML_toy/output/txtFileWithOutput.txt

Summary text output: /Users/zg/CML_toy/output/txtFileWithOutput_Summary.txt

BACKGROUND

ESTIMATION OF EFFECTS

The MCF is a comprehensive causal machine learning estimator for the estimation of treatment effects at various levels of granularity, from the average effect at the population level to very fine grained effects at the (almost) individual level. Since effects at the higher levels are obtained from lower level effects, all effects are internally consistent. Recently, the basic package has been appended for new average effects as well as for an optimal policy module. Effect estimation is implemented for identification by unconfoundedness as well as by instrumental variables. While unconfoundedness estimation can deal with multiple treatments, instrumental variable estimation is restricted to binary instruments and binary treatments. The basis of the MCF estimator is the causal forest suggested by Wager and Athey (2018). Their estimator has been changed in several dimensions which are described in Lechner (2018). The main changes relate to the objective function as well as to the aggregation of effects. Lechner and Mareckova (2024) provide the asymptotic guarantees for the MCF and compare the MCF, using a large simulation study, to competing approaches like the Generalized Random Forest (GRF, Athey, Tibshirani, Wager, 2019) and Double Machine Learning (DML, Chernozhukov, Chetverikov, Demirer, Duflo, Hansen, Newey, Robins, 2018, Knaus, 2022). In this comparison the MCF faired very well, in particular, but not only, for heterogeneity estimation. Some operational issues of the MCF are discussed in Bodory, Busshof, Lechner (2022). There are several empirical studies using the MCF, like Cockx, Lechner, Boolens (2023), for example.

References

- Athey, S., J. Tibshirani, S. Wager (2019): Generalized Random Forests, The Annals of Statistics, 47, 1148-1178.
- Athey, S., S. Wager (2019): Estimating Treatment Effects with Causal Forests: An Application, Observational Studies, 5, 21-35.
- Bodory, H., H. Busshoff, M. Lechner (2023): High Resolution Treatment Effects Estimation: Uncovering Effect Heterogeneities with the Modified Causal Forest, Entropy, 24, 1039.
- Bodory, H., F. Mascolo, M. Lechner (2024): Enabling Decision Making with the Modified Causal Forest: Policy Trees for Treatment Assignment, Algorithm, 17, 318.
- Chernozhukov, V., D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, J. Robins (2018): Double/debiased machine learning for treatment and structural parameters, Econometrics Journal, 21, C1-C68.

Modified Causal Forest: Estimation

- Cockx, B., M. Lechner, J. Bollens (2023): Priority to unemployed immigrants? A causal machine learning evaluation of training in Belgium, Labour Economics, 80, Article 102306.
- Knaus, M. (2022): Double Machine Learning based Program Evaluation under Unconfoundedness, Econometrics Journal.
- Lechner, M. (2018): Modified Causal Forests for Estimating Heterogeneous Causal Effects, arXiv.
- Lechner, M. (2023): Causal Machine Learning and its Use for Public Policy, Swiss Journal of Economics & Statistics, 159:8.
- Lechner, M., J. Mareckova (2024): Comprehensive Causal Machine Learning, arXiv.
- Lechner, M., J. Mareckova (2025): Comprehensive Causal Machine Learning with Instrumental Variables, mimeo.
- Wager, S., S. Athey (2018): Estimation and Inference of Heterogeneous Treatment Effects using Random Forests, Journal of the American Statistical Association, 113:523, 1228-1242.

Modified Causal Forest: Estimation

Section 2: MCF estimation

METHOD

Standard MCF method used. Nearest neighbour matching performed using the Prognostic Score.
Feature selection not is used.
Local centering is used.
Common support is enforced.

VARIABLES

Outcome: sal_9

Treatment: ptype (with values 0 1 2)

Ordered confounders: sex, specia_cw, age, school, voc_deg, reg_al, reg_ser, reg_pro, reg_agri, sect_al, prof_al, unem_x0, olf_x0, empl_x0, earn_x0, emplx1_1, emplx1_2, emplx1_3, emplx1_4, emplx2_1, emplx2_2, emplx2_3, emplx2_4, earnx1_1, earnx1_2, earnx1_3, earnx1_4, earnx2_1, earnx2_2, earnx2_3, earnx2_4, Imp_cw

Unordered (categorical) confounders: nation, region

EFFECTS ESTIMATED

Average Treatment Effect (ATE), Individualized Average Treatment Effect (IATE)

NOTE on unordered variables:

One-hot-encoding (dummy variables) is not used as it is expected to perform poorly with trees: It may lead to splits of one category versus all other categories. Instead the approach used is analogous to the one discussed in Chapter 9.2.4 of Hastie, Tibshirani, Friedman (2013), The Elements of Statistical Learning, 2nd edition.

Section 2.1: MCF Training

Training uses 10 CPU cores.

Section 2.1.1: Preparation of training data (mcf training)

METHOD

Variables without variation are removed.
Variables that are perfectly correlated with other variables are removed.
Dummy variables with less than 10 observations in the smaller group are removed.
Rows with any missing values for variables needed for training are removed.

RESULTS

No relevant variables were removed.
Sample size of training data: 2899 (no observations removed).

Modified Causal Forest: Estimation

Section 2.1.2: Common support (mcf training)

METHOD

The common support analysis is based on checking the overlap in the out-of-sample predictions of the propensity scores (PS) for the different treatment arms. PSs are estimated by random forest classifiers. Overlap is operationalized by computing cut-offs probabilities of the PSs (ignoring the first treatment arm, because probabilities add to 1 over all treatment arms). These cut-offs are subsequently also applied to the data used for predicting the effects.

Overlap is determined by the min / max rule.

Cut-offs for PS are widened by 0.05.

Out-of-sample predictions are generated by 5-fold cross-validation.

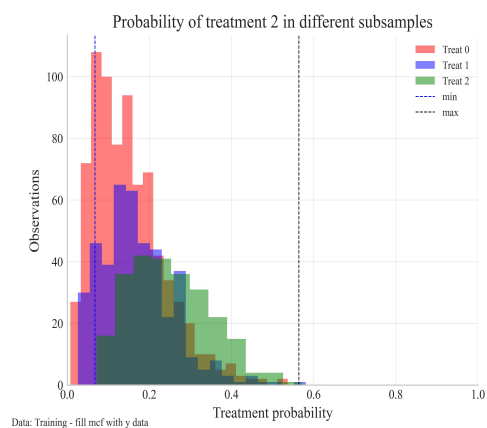
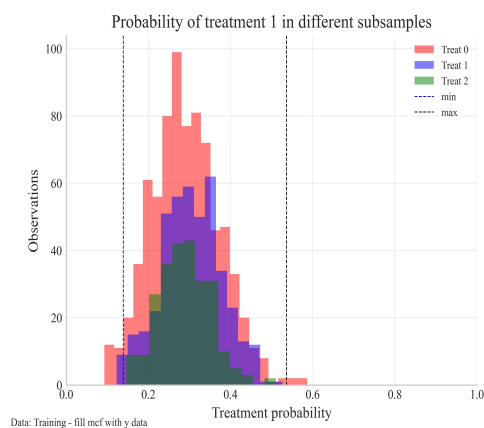
RESULTS

Share of observations deleted: 13.14%

Number of observations remaining: 2518

WARNING: Check output files whether the distribution of the features changed due to the deletion of part of the data.

Common support plots



Section 2.1.3: Local centering (mcf training)

METHOD

Local centering is based on training a regression to predict the outcome variable conditional on the features (without the treatment). The regression method is selected among various versions of Random Forests, Support Vector Machines, Boosting methods, and Neural Networks of scikit-learn. The best method is selected by minimizing their out-of-sample Mean Squared Error using 5-fold

Modified Causal Forest: Estimation

cross-validation. The full set of results of the method selection step are contained in /Users/zg/CML_toy/output/txtFileWithOutput.txt.

The respective out-of-sample predictions are subtracted from the observed outcome in the training data used to build the forest. These out-of-sample predictions are generated by 5-fold cross-validation.

RESULTS

Out-of-sample fit for Random Forest of $E_{y|x}$ (R^2) for sal_9: 42.99%

Section 2.1.4: Forest

METHOD and tuning parameters

Method used for forest building is MSE & MCE Penalty "MSE of treatment variable".

The causal forest consists of 1000 trees.

The minimum leaf size is 5.

The number of variables considered for each split is 12.

The share of data used in the subsamples for forest building is 67%.

The share of the data used in the subsamples for forest evaluation (outcomes) is 100%.

Alpha regularity is set to 10%.

sal_9_lc is the outcome variable used for splitting (locally centered).

The features used for splitting are sex specia_cw age school voc_deg reg_al reg_ser reg_pro reg_agri sect_al prof_al unem_x0 olf_x0 empl_x0 earn_x0 emplx1_1 emplx1_2 emplx1_3 emplx1_4 emplx2_1 emplx2_2 emplx2_3 emplx2_4 earnx1_1 earnx1_2 earnx1_3 earnx1_4 earnx2_1 earnx2_2 earnx2_3 earnx2_4 Imp_cw nation region.

RESULTS

Each tree has on average 109.81 leaves.

Each leaf contains on average 7.6 observations. The median # of observations per leaf is 7.

The smallest leaves have 5 observations.

The largest leaf has 61 observations.

23.96% of the leaves were merged when populating the forest with outcomes from the honesty sample.

Modified Causal Forest: Estimation

Section 2: MCF estimation

Section 2.2: MCF Prediction of Effects

Training uses 10 CPU cores.

Section 2.2.1: Common support (mcf prediction)

Share of observations deleted: 13.97%

Number of observations remaining: 2495

WARNING: Check output files whether the distribution of the features changed due to the deletion of part of the data.

Section 2.2.2: Results

GENERAL REMARKS

The following results for the different parameters are all based on the same causal forests (CF). The combination of the CF with the potentially new data provided leads to weight matrices. These matrices may be large requiring some computational optimisations, such as processing them in batches and saving them in a sparse matrix format. One advantage of this approach is that aggregated effects (ATE, GATE, BGATE) can be computed by aggregation of the weights used for the IATE. Thus a high internal consistency is preserved in the sense that IATEs will aggregate to GATEs, which in turn will aggregate to ATEs.

ESTIMATION

Weights of individual training observations are truncated at 5.00%. Aggregation of IATEs to ATE and GATEs may not be exact due to weight truncation.

INFERENCE

Inference is based on using the weight matrix. Nonparametric regressions are based on k-nearest neighbours.

NOTE

Treatment effects for specific treatment groups (so-called treatment effects on the treated or non-treated) can only be provided if the data provided for prediction contains a treatment variable (which is not required for the other effects).

Section 2.2.2.1: ATE

Modified Causal Forest: Estimation

RESULT

ATE for sal_9

Comparison	Effect	SE	t-value	p-value (%)	Sig.
1 vs 0	8588.516	1053.027	8.16	0.0	****
2 vs 0	678.859	2080.924	0.33	74.14	
2 vs 1	-7909.657	2306.565	3.43	0.06	****

Note: *, **, ***, **** denote significance at the 10%, 5%, 1%, 0.1% level. The results for the potential outcomes can be found in the output files.

Section 2.2.2.2: IATE

This section contains parts of the descriptive analysis of the IATEs. Use the analyse method to obtain more descriptives of the IATEs, like their distribution, and their relations to the features.

RESULTS

Outcome variable: sal_9

Comparison	Mean	Median	Std	Effect > 0	mean(SE)	sig 10%	sig 5%	sig 1%
1 vs 0	8589.70021	9909.29436	6548.05961	86.77%	2142.27384	83.73%	79.40%	72.22%
2 vs 0	663.28149	928.84059	6594.47404	54.99%	3979.18861	38.08%	27.78%	10.38%
2 vs 1	-7926.41873	-5466.87882	8112.41508	14.11%	4479.19527	41.68%	36.71%	31.34%