

Supplementary Material 1 for the article:
**Improved outcome prediction across data sources
through robust parameter tuning**

Nicole Schüller^{*,1}, Anne-Laure Boulesteix¹, Bernd Bischl², Kristian Unger^{3,4},
Roman Hornung¹

¹ Institute for Medical Information Processing, Biometry and Epidemiology, University of Munich,
Marchioninstr. 15, 81377 Munich, Germany

² Department of Statistics, University of Munich, Ludwigstrasse 33, 80539 Munich, Germany

³ Research Unit Radiation Cytogenetics, Helmholtz Zentrum Munich, German Research Center for
Environmental Health GmbH, Neuherberg, Germany

⁴ Department of Radiation Oncology, University Hospital, University of Munich, Munich, Germany

*To whom correspondence should be addressed: nschueller@ibe.med.uni-muenchen.de

A Overview shown in the main paper extended by the procedures that include the external data set for training

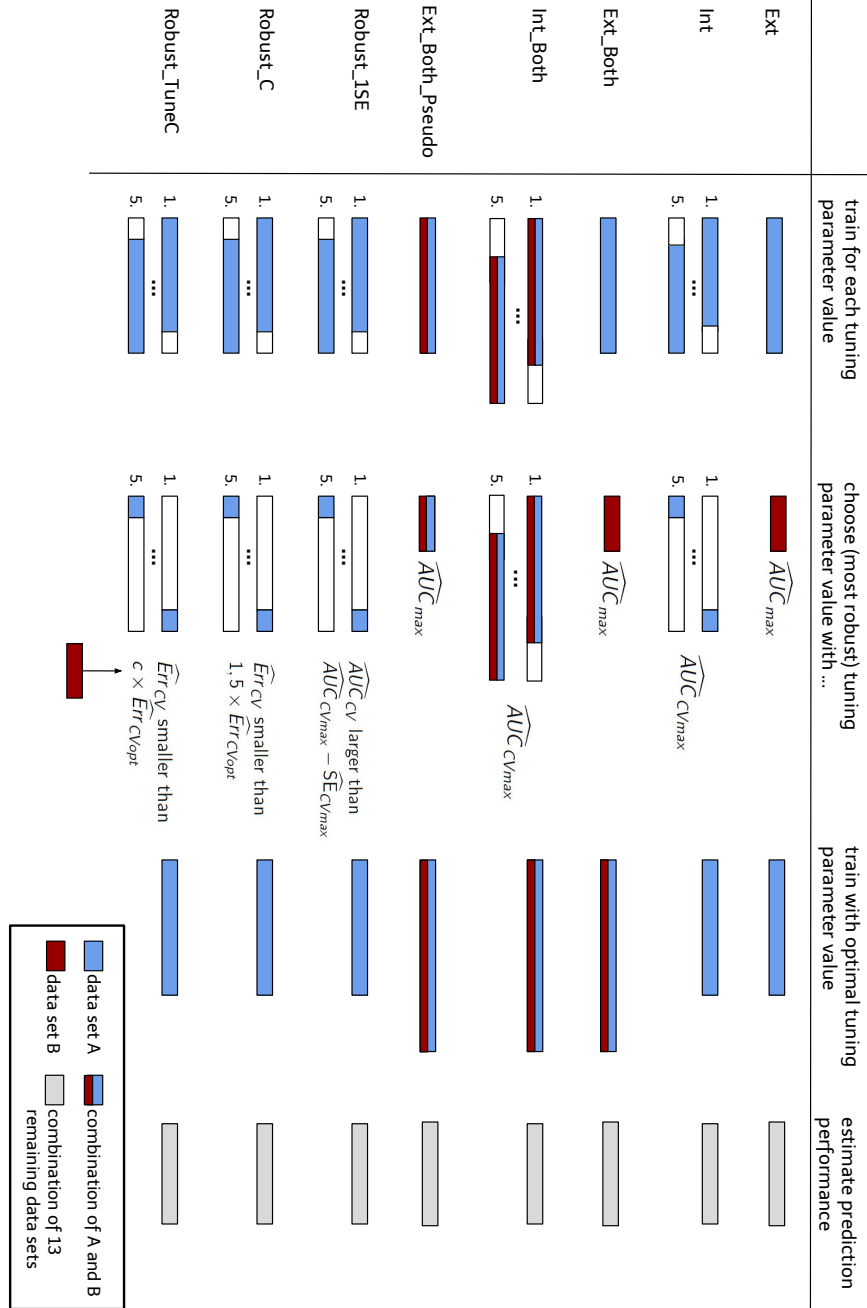


Fig. S1: Overview of the practically motivated approaches to external / internal tuning and the procedures for robust tuning including the procedures that include the external data set for training. Note that the two procedures Robust_C and Robust_TuneC are depicted in a simplified way. See sections 2.5.3 and 2.5.4 of the main paper for details.

B Procedures using both A and B for training

The procedures `Ext` and `Int` described in section 2.5.1 of the main paper use the external data set only for choosing the optimal tuning parameter value (`Ext`) or not at all (`Int`). It might, however, be worthwhile to use data set B also for training. We included two such procedures in our comparison study.

In the first one, `Ext_Both`, external tuning is performed (as in `Ext`). Then A and B are combined with batch effects adjustment using ComBat (Johnson et al., 2007) to finally train the prediction rule with the optimal tuning parameter value.

In the second variant, `Int_Both`, A and B are first combined (again performing ComBat to adjust for batch effects). Then internal tuning is applied to the combined data (i.e., `Int` is applied to the combination of A and B). This method is followed by researchers who are only aware of internal tuning and want to use all observations (from A and B) for training the prediction rule.

Finally, we also include an additional variant, `Ext_Both_Pseudo`, of the first approach `Ext_Both`. This procedure would not be used in practice. Its aim is to assess whether it is important that, in `Ext_Both`, the external data set B comes from a different distribution than A. To assess this, we proceed as follows. As with `Int_Both` we first combine A and B adjusting for batch effects using ComBat. Subsequently, we randomly split the combined data into two parts. The first part has the same size n_A as data set A and the second part has the same size n_B as data set B. Subsequently, the `Ext_Both` procedure is applied to these two parts, which play the role of training and (pseudo) external tuning data set, respectively. In order to reduce the dependency of the results of `Ext_Both_Pseudo` on the specific random splitting, we again repeat the procedure for 10 random splits. The 10 optimized tuning parameter values and the 10 obtained AUC values are subsequently averaged. This approach is essentially the same as `Ext_Both`, except that the training and (pseudo) external tuning data set follow the same (mixture) distribution.

C Extended results of the conceptual comparison of external and internal tuning

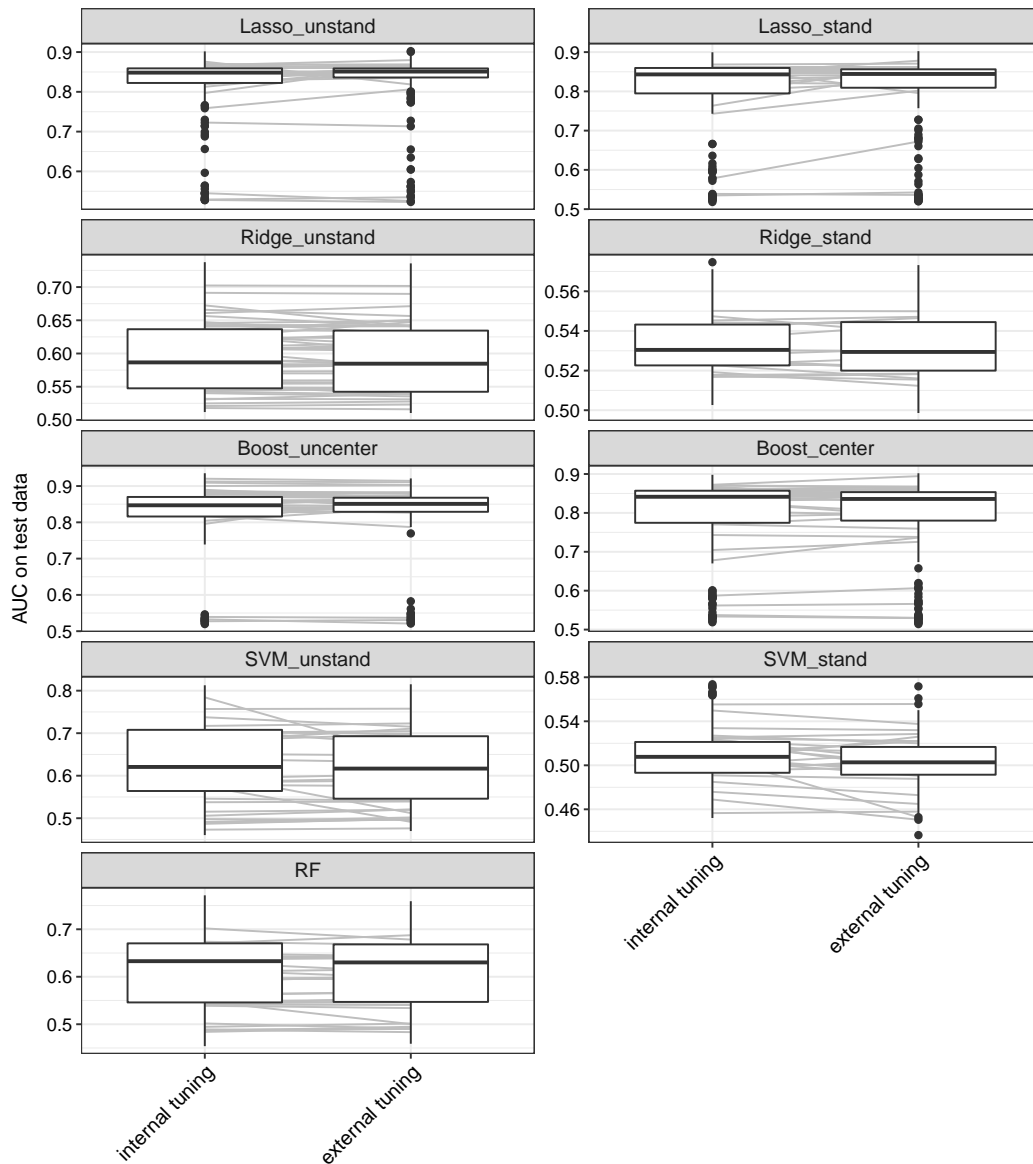


Fig. S2: Extended results: test data estimated prediction performance estimates in the conceptual comparison of external and internal tuning. The grey lines connect the values of pairs that share the same training data sets, where in each case, for the sake of clarity, we do not show a line for each of the pairs, but merely for a random subset of 30 pairs.

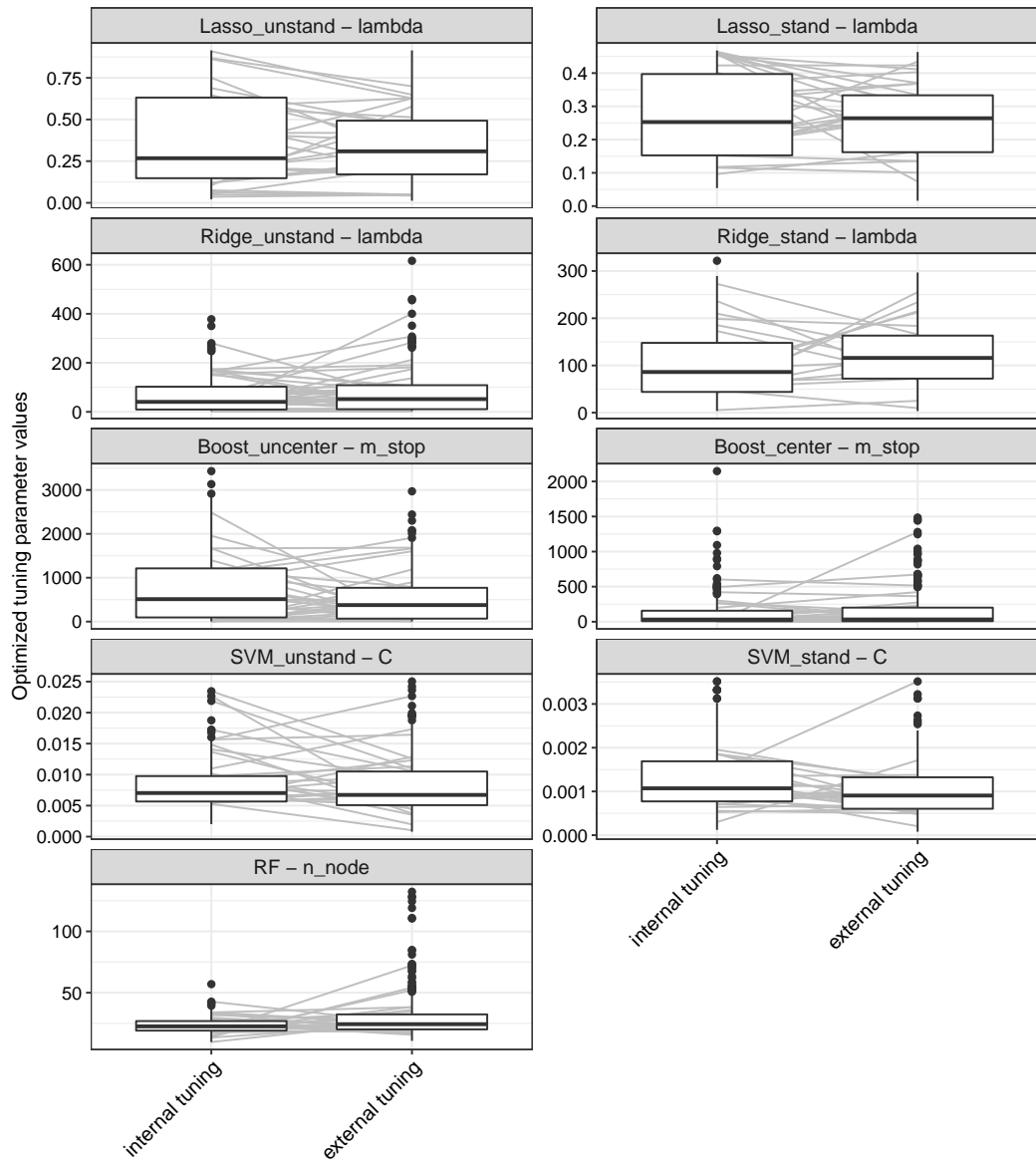


Fig. S3: Extended results: chosen tuning parameter values in the conceptual comparison of external and internal tuning. The grey lines connect the values of pairs that share the same training data sets, where in each case, for the sake of clarity, we do not show a line for each of the pairs, but merely for a random subset of 30 pairs.

D Extended results: prediction performance estimates in the comparison of the practically motivated tuning approaches

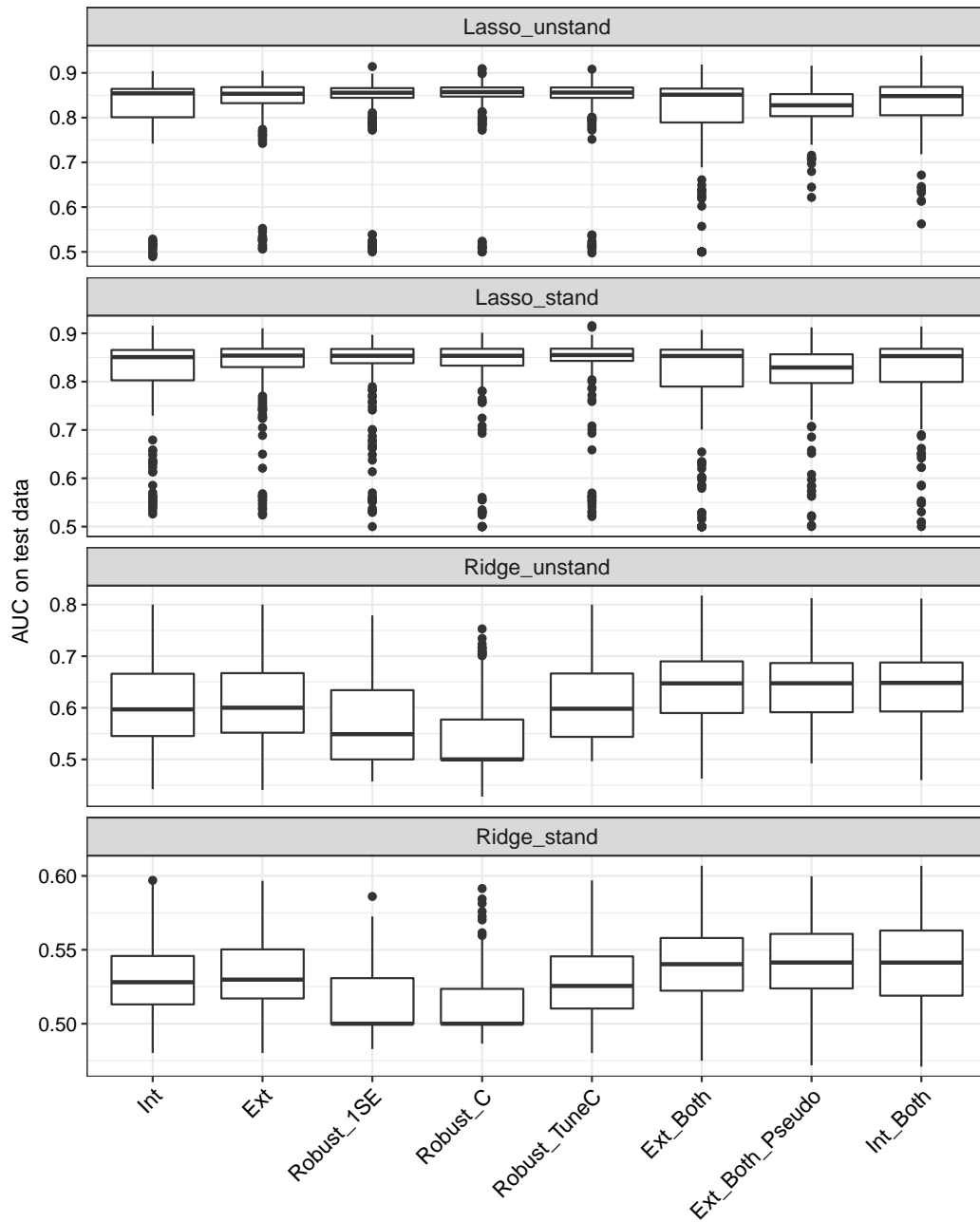


Fig. S4: Extended results: test data estimated prediction performance estimates in the comparison of the practically motivated tuning approaches - I

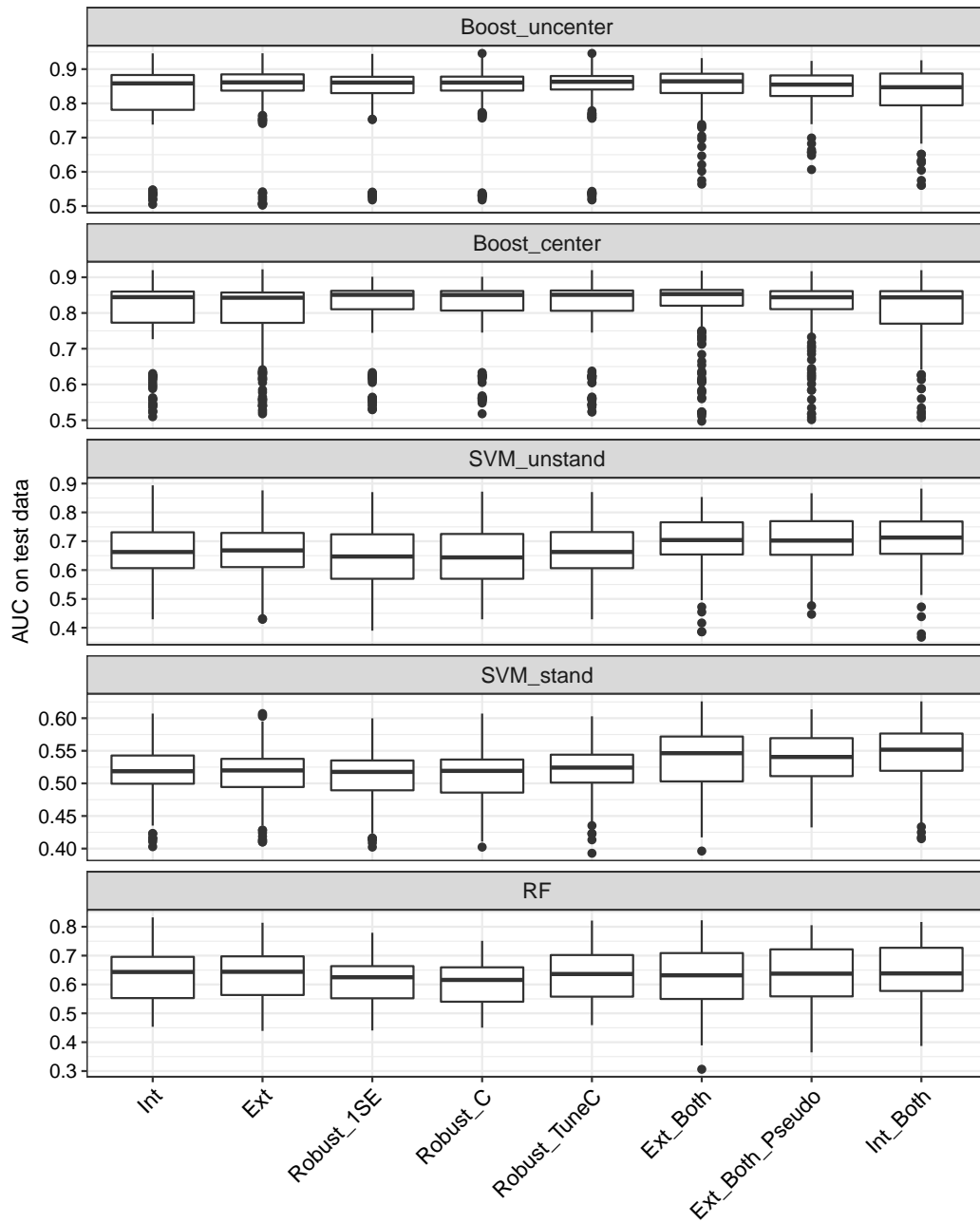


Fig. S5: Extended results: test data estimated prediction performance estimates in the comparison of various practically motivated tuning approaches - II

E Extended results: chosen tuning parameter values in the comparison of the practically motivated tuning approaches

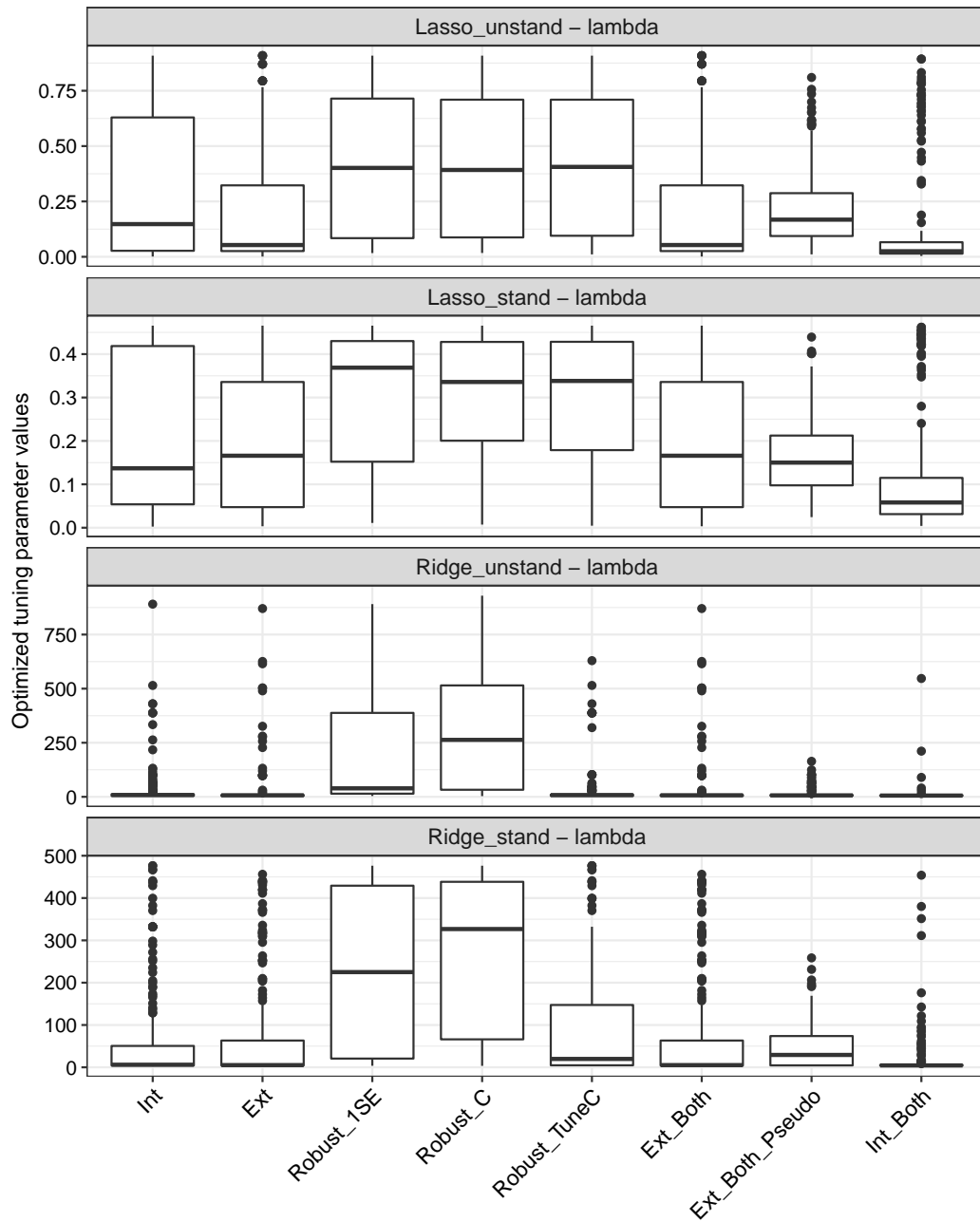


Fig. S6: Extended results: chosen tuning parameter values in the comparison of the practically motivated tuning approaches - I

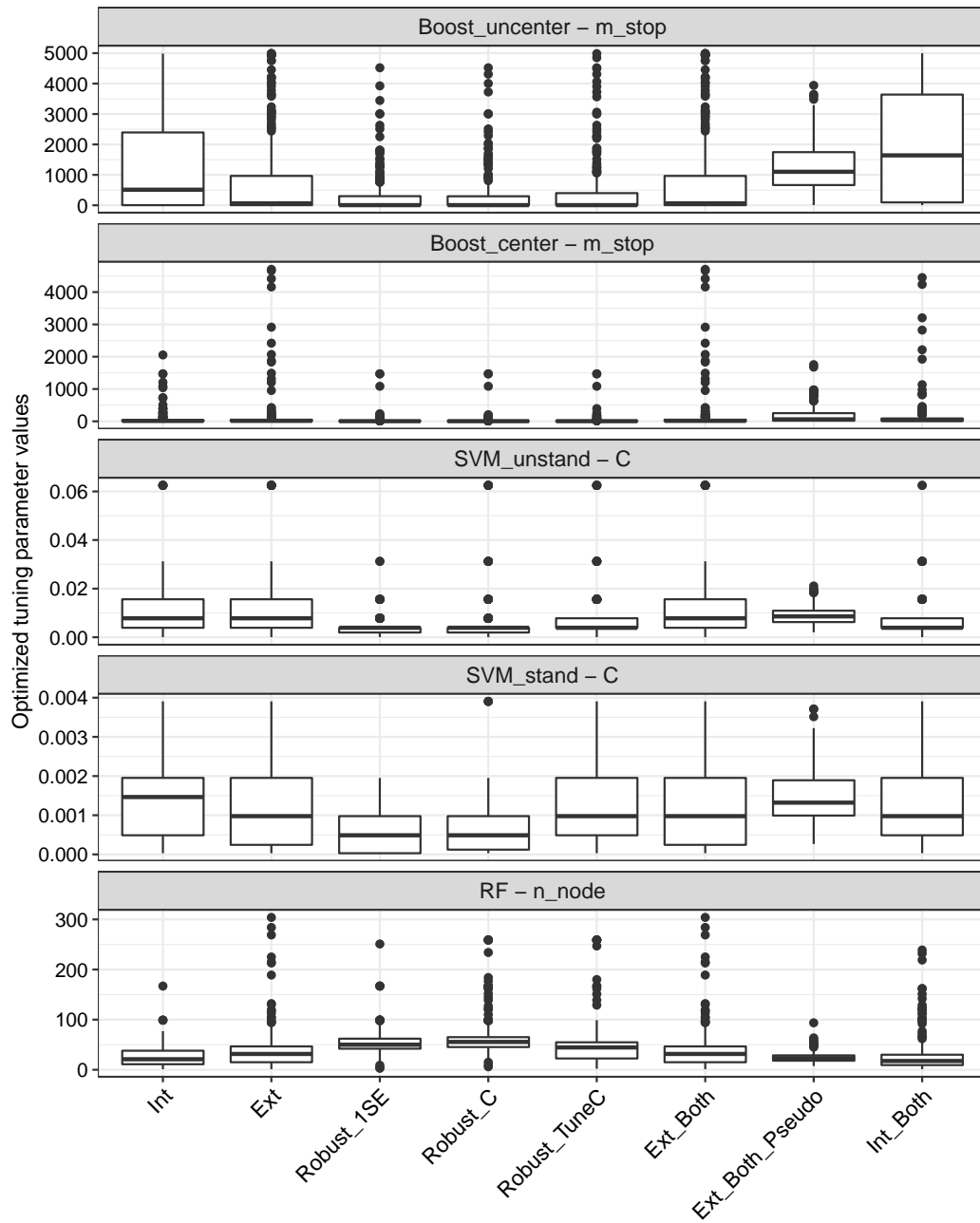


Fig. S7: Extended results: chosen tuning parameter values in the comparison of the practically motivated tuning approaches - II

F Robust_TuneC: Chosen c values

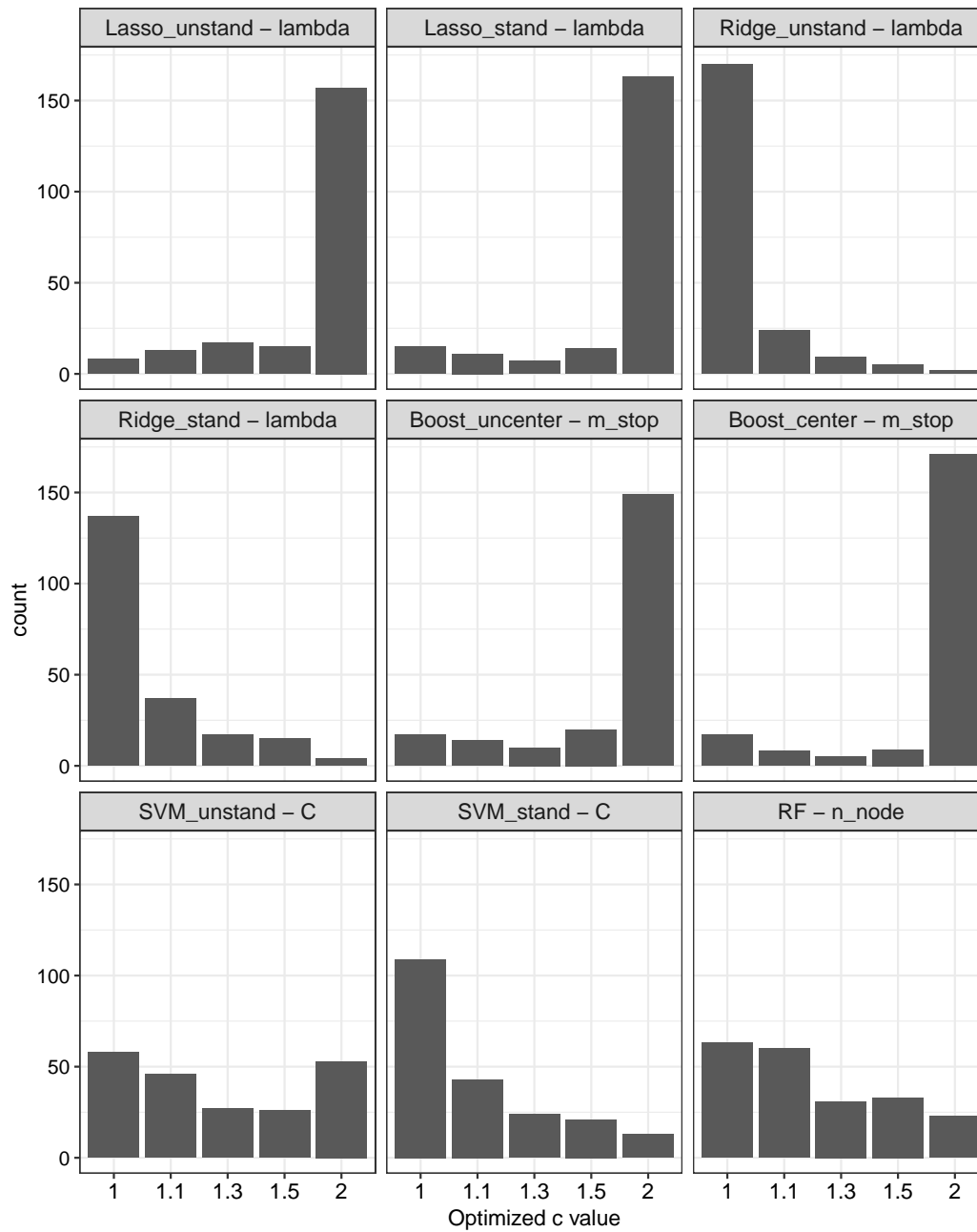


Fig. S8: Frequencies of c values chosen from the grid used for Robust_TuneC in the comparison of the practically motivated tuning approaches

G Discussion of the results obtained with the procedures that include the external data set for training

The three approaches that include the external data set for training, `Ext_Both`, `Ext_Both_Pseudo`, and `Int_Both` (see section B), perform slightly better than the other approaches for some of the prediction methods (see section D). More precisely, in cases in which the cross-study prediction performance is strong overall, these approaches do not perform better than competitors. `Ext_Both` performs slightly better than `Ext_Both_Pseudo` for some prediction methods, while for others there are no notable differences. The latter suggests that it is not very beneficial to perform external tuning to choose the tuning parameter value and subsequently combine the training data set with the external data set for training the prediction rule. For some prediction methods `Ext_Both` outperforms `Int_Both`, while there are no systematic differences for the remaining prediction methods. `Ext_Both_Pseudo` seems to perform better than `Int_Both` for `Boost_uncenter` and `Boost_center`. This is, however, likely due to the fact that the variance of the AUC estimates is reduced for `Ext_Both_Pseudo`, as in the case of this approach the AUC estimates do not represent AUC values obtained from single evaluations of the test data. Instead, they are obtained by averaging the AUC values obtained from ten iterations, where each of these corresponds to a different split into “pseudo training data set” and “pseudo external data set”, see again section B.

H Extended results of the additional study: optimistic bias by using the external data set for both tuning and prediction performance estimation

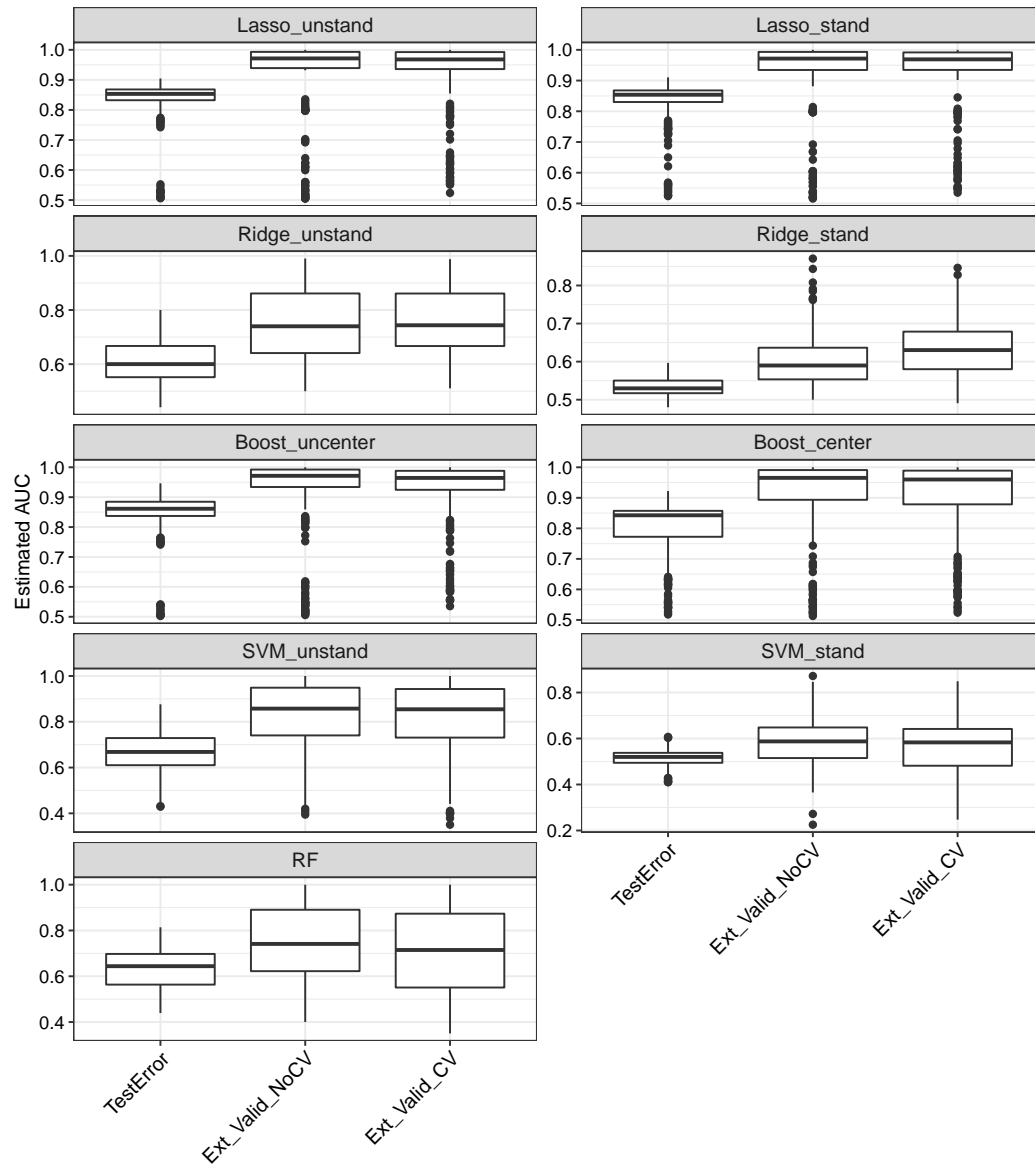


Fig. S9: Extended results: comparison of test data estimated prediction performance estimates with those obtained by prediction performance estimation approaches that use the external data set both for choosing the tuning parameter value and for prediction performance estimation

References

Johnson WE, Rabinovic A, Li C (2007) Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* 8:118–127