
A Large-Scale Neutral Comparison Study of Survival Models on Low-Dimensional Data

Lukas Burk^{1,2,3,4} John Zobolas⁵ Bernd Bischl^{1,4}
 Andreas Bender^{1,4} Marvin N. Wright^{2,3,6} Raphael Sonabend^{7,8}

¹ Department of Statistics, LMU Munich, Munich, Germany

² Leibniz Institute for Prevention Research and Epidemiology - BIPS, Bremen, Germany

³ Faculty of Mathematics and Computer Science, University of Bremen, Bremen, Germany

⁴ Munich Center for Machine Learning, Munich, Germany

⁵ Department of Cancer Genetics, Institute for Cancer Research, Oslo University Hospital, Oslo, Norway

⁶ Department of Public Health, University of Copenhagen, Copenhagen, Denmark

⁷ OSPO Now, London, UK

⁸ Imperial College London, London, UK

{lukas.burk,bernd.bischl,andreas.bender}@stat.uni-muenchen.de
 ioannisz@uio.no
 wright@leibniz-bips.de
 raphaelsonabend@gmail.com

Abstract

This work presents the first large-scale neutral benchmark experiment focused on single-event, right-censored, low-dimensional survival data. Benchmark experiments are essential in methodological research to scientifically compare new and existing model classes through proper empirical evaluation. Existing benchmarks in the survival literature are often narrow in scope, focusing, for example, on high-dimensional data. Additionally, they may lack appropriate tuning or evaluation procedures, or are qualitative reviews, rather than quantitative comparisons. This comprehensive study aims to fill the gap by neutrally evaluating a broad range of methods and providing generalizable conclusions. We benchmark 18 models, ranging from classical statistical approaches to many common machine learning methods, on 32 publicly available datasets. The benchmark tunes for both a discrimination measure and a proper scoring rule to assess performance in different settings. Evaluating on 8 survival metrics, we assess discrimination, calibration, and overall predictive performance of the tested models. Using discrimination measures, we find that no method significantly outperforms the Cox model. However, (tuned) Accelerated Failure Time models were able to achieve significantly better results with respect to overall predictive performance as measured by the right-censored log-likelihood. Machine learning methods that performed comparably well include Oblique Random Survival Forests under discrimination, and Cox-based likelihood-boosting under overall predictive performance. We conclude that for predictive purposes in the standard survival analysis setting of low-dimensional, right-censored data, the Cox Proportional Hazards model remains a simple and robust method, sufficient for practitioners.

1 Introduction

Survival analysis is an important branch of statistics for data where the outcome is the time until an event of interest occurs. Such data often exhibits incomplete information about the outcome, for example due to censoring. Traditionally applied in medical research to estimate how patient survival

relates to features, it has been applied in a broad range of applications across various domains. By effectively incorporating information from both completed and ongoing (censored) cases, survival analysis can yield accurate and informative predictions. This capability is invaluable in fields such as medicine, finance, and in different industrial sectors, where risk prediction is an important decision making component. Many methods have been introduced in this field, from the Cox Proportional Hazards (CPH) [26] and Accelerated Failure Time (AFT) model [51] to tree-based methods including Random Survival Forests (RSFs) [48] and Gradient Boosting Machines (GBMs) [33, 25] as well as many others.

Throughout this paper, we consider only the right-censored survival setting. This restriction is due to a lack of comprehensive support for other settings (e.g. left censoring or competing risks) in available methods and especially their software implementation. In a nutshell, right-censoring occurs when some subjects did not experience the event of interest either because of drop-out or end of study (both assumed to be unrelated to the event of interest here). Formally, let $Y_i \sim F_Y; Y_i > 0; i = 1, \dots, n$ be a random variable representing the event times and $C_i \sim F_C; C_i > 0$ the censoring time. In right censored data, we do not observe realizations of Y_i but rather of the tuple $(T_i = \min(Y_i, C_i), D_i = \mathbb{I}(Y_i \leq C_i))$. The goal of survival analysis is to obtain estimates/predictions for the distribution F_Y , or quantities derived from it, e.g. $\mathbb{E}_Y(Y)$, based on realizations of (T_i, D_i) ; see Section 3.3 for definition of prediction types. The observed data is given by tuples (t_i, d_i, \mathbf{x}_i) , $i = 1, \dots, n$, where t_i is the *observed* outcome time (either event or censoring time, whichever occurred first), d_i is the status indicator (0 if observation is censored, 1 if the event of interest was observed) and \mathbf{x}_i is the feature vector.

Contributions: This paper introduces what we believe to be the first large-scale comparison study for single-event, right-censored survival data due to the number and range of datasets (32), models (18), tuning measures (2), and evaluation measures (8) included. We benchmark survival techniques in the the low-dimensional setting, which represents a type of data that practitioners often encounter. While ex ante one might suspect that machine learning methods will not perform well in low-dimensional settings, there are no established rules for sample sizes required by machine learning methods in the context of survival analysis. The size of our study and the inclusion of many diverse, and hopefully representative, data sets should ensure that its results can be generalized within other right-censored, low-dimensional settings [43]. Furthermore, we run a “neutral” benchmark in accordance with the guidelines laid out by Boulesteix *et al.* [16] which we detail in Section 3.1. Based on our review of the literature (see Section 2), there is no other study: a) with a comparable number of datasets or methods; b) that compared methods after sufficient tuning for both discrimination and overall predictive ability (as measured by scoring rules); or c) that neutrally compares methods. Finally, our collection of datasets is available as an OpenML benchmark suite [99] and our hyperparameter search spaces will be available in a forthcoming release of *mlr3tuningspaces* [9].

2 Literature Review

The experiments described in this paper provide a comparison of both classical and machine learning (ML) survival models in a low-dimensional setting. We use these broad terms for model classes in analogy to the taxonomy provided by [101], where the term “classical” refers to semi- and fully parametric methods such as the CPH and AFT models or derivatives thereof, including penalized variants. “ML” here refers to non-linear and non-parametric methods ranging from tree-based methods including RSFs, boosting approaches such as GBMs or likelihood boosting (CoxBoost), to Survival Support Vector Machines (SSVMs), Artificial Neural Networks (ANNs) and Deep Learning (DL) methods. Although there is no objective and generally accepted definition to determine whether a dataset is “high-dimensional”, we colloquially define it to refer to scenarios where the number of features exceeds the number of observations ($p > n$).

Historically, surveys, reviews, and analytical comparisons of survival models can be grouped into: i) empirical comparisons of models with limited scope; and ii) qualitative surveys without benchmark experiments.

We provide a short overview of the literature for comparisons of survival models, with an extended review available in Appendix A. Papers that empirically compare survival models are further separated into studies that: i) compare ‘classical’ models only; ii) compare multiple ML and classical model

classes; iii) compare one novel model (or class) to one or more baseline models; and iv) exclusively focus on high-dimensional data.

Comparisons of Classical Models often compare CPH and AFT models, including [74, 35, 106, 40] showing both methods yielding similar hazard ratios but without evaluation metrics on independent test data, relying on graphical procedures to draw conclusions; [28] additionally compare flexible Cox models including splines using time-dependent AUC, MSE, and MAE, finding that the CPH with penalized splines outperformed the other models.

Comparisons of ML and Classical Models The experiments carried out in this paper belong in this category. Only two prior experiments could be found that neutrally benchmarked more than one ML model class on low-dimensional data. Kattan (2003) [53] benchmarked tree-based models, ANNs and CPH with Harrell’s C-index across three datasets. The models are compared for significant differences using 50 times repeated nested cross-validation but do not clarify their tuning procedure. Boxplots across all replications indicate that no machine learning model outperformed the CPH. The authors note the small number of datasets used for comparison as their primary limitation. Zhang *et al.* (2021) [107] compare classical and ML methods, taking into account feasibility and computational efficiency for various tasks in the biomedical field. Methods are evaluated on six clinical and 16 omics datasets using 11 metrics, including time-dependent AUC, Brier score and multiple variations of the C-index. Methods were applied with specific hyperparameter settings without tuning, which limits the generalizability of their results.

Comparisons of a Novel Model Class include [71, 76] benchmarking newly developed ANNs against CPH and [37] comparing new SVMs against CPH, where neither found significant differences. [49] compare a novel implementation of oblique RSFs (“aorsf”) to the previous implementation, as well as other RSFs, GBMs, penalized CPH, and ANNs. Using Harrell’s C and ISBS, they find aorsf to outperforming GBMs and penalized CPH but with only minimal tuning applied for some models.

Comparisons on High-Dimensional Data gained popularity in the survival literature, with many recent studies focusing on the area of multi-omics data [108]. [43] perform a large-scale benchmark including penalized regression, GBMs, and RSFs, evaluating on Harrell’s C and ISBS but not finding any significant differences. [92] compare a similar group of models evaluating on Harrell’s C only, finding few significant differences. [104] compare DL methods, RSFs and CPH, evaluating on Antolini’s C-index and ISBS and noting a lack of noise-resistance of all models.

3 Benchmark Experiments

3.1 Study Design

The experiments in this study are designed to assess the status quo of survival models, including both classical and machine learning approaches. In order to achieve this objective, this study aims to be a “neutral comparison study” [16]. Following the guidelines put forward by Boulesteix *et al.*:

Focus on model comparison: The focus of this study is on model comparison rather than on the examination of a novel model. We do not favor one dataset over another and draw conclusions across all datasets instead of trying to find data sets in which the models performed well.

The authors are neutral: At least one representative of all methods compared in this experiment was contacted, and hyperparameter configurations were discussed with all who responded. Every party with a personal interest in the results was provided equal opportunity to influence the experiment and thereby making sure that there was no bias involved in model configuration. We are grateful for the maintainers’ time supporting this effort.

Model, performance measures and data are chosen in a rational way: The study is designed to assess the status quo, which excludes models and measures that have been published without proper peer-review. We use one primary measure each for discrimination and overall predictive ability, while additional measures are reported for comparison purposes. In order to be in line with common procedures and to allow for general comparability, we also assess models by measures even if these are known to be flawed, e.g., increasing bias of Harrell’s C for increasing censoring percentages (see, for example, [83]). The inclusion criteria for datasets were as follows: We use real-world datasets that include at least two features, a right-censoring indicator and a survival time, have at least 100 observed events, and which do not qualify as high-dimensional, i.e. have fewer

features than observations. We explicitly exclude datasets with competing risk endpoints, recurrent events, or other non-standard settings such as truncation or left-censoring. No quota was specified regarding censoring proportions in the datasets.

Implementation, Reproducibility, and Accessibility Experiments were conducted on R 4.2.2 on the Beartooth Computing Environment [8]. All code required to run the experiments and generate the results is available in a public GitHub repository¹ licensed under GPL-3. Further details on software used are available in Appendix E. For additional reproducibility, our hyperparameter search spaces will be published with the next release of *mlr3tuningspaces* [9] and our datasets will be available as an OpenML benchmark suite [99].

3.2 Models and Configurations

The models compared in this experiment were chosen by identifying commonly used models with readily available implementations: i) Kaplan-Meier (KM) [52]; ii) Nelson-Aalen (NA) [1]; iii) Akritas Estimator (AK) [2]; iv) Cox PH (CPH) [26]; v) CV Regularized CPH (GLMN) [84]; vi) Penalized (Pen) [36]; vii) Parametric AFT (AFT) [51]; viii) Flexible Splines (Flex) [82]; ix) Random Survival Forest (RFSRC) [48]; x) Random Survival Forest (RAN) [48, 105]; xi) Conditional Inference Forest (CIF) [46]; xii) Oblique Random Survival Forest (ORSF) [49]; xiii) Relative Risk Tree (RRT) [17]; xiv) Model-Based Boosting (MBST) [20]; xv) CoxBoost (CoxB) [12]; xvi) XGBoost with Cox objective (XGBCox) [25]; xvii) XGBoost with AFT objective (XGBAFT) [7]; and xviii) SSVM-Hybrid (SSVM) [97].

The full table of all models including respective software packages and versions is given in Appendix C. In our selection we focused on well-established models with robust implementations, provided as well maintained packages or wrapper functions within benchmarking software. This excludes some recently proposed DL based methods like DeepSurv [55] and DeepHit [68], which have higher computational complexity, require intensive tuning, and in initial experiments could not be evaluated reliably within our benchmark suite. The KM and NA estimators are used as non-parametric baselines while AK acts as a more flexible baseline as it estimates a conditional survival function without assuming uninformative censoring. An additional table in Appendix D lists the models' hyperparameter and pre-processing configurations. Some models could be considered equivalent as they implement the same method, e.g. RAN and RFSRC both implement Random Survival Forests, yet for the purposes of this benchmark we consider the implementation a part of the model comparison.

3.3 Resampling, Tuning, Prediction Types, and Pre-Processing

Resampling is performed as nested cross-validation with five outer and three inner resampling folds for unbiased generalization error estimates [14]. Stratified resampling is used to preserve the proportion of censoring in respective folds. The number of outer resampling folds is reduced to four if necessary to prevent folds with less than 30 observations, which affected only the veteran dataset from the survival package.

Tuning is performed on the inner resampling folds of the nested cross-validation. We use random search with 50ρ iterations where ρ is the dimensionality of the search space (ranging from 1 to 8). For example, one parameter of GLMN is tuned for 50 iterations, whereas six parameters of ORSF are tuned for 300 iterations. This method grants each method the opportunity to be tuned to the same relative amount, i.e. 50 iterations per tunable parameter. When the tuning space was finite and smaller than 50, we used exhaustive grid search to achieve the same tuning result but with lower computational cost; this applied to learners AFT (tuning across one of three distribution families), Flex (tuning $k \in 1 \dots 10$, and RRT (tuning `minbucket` $\in 5 \dots 50$). The secondary stopping criterion for the tuning process was a time limit of 150 hours, which ensured that one outer resampling iteration (tuning and final model fitting), could complete within seven days, a constraint imposed by the computational environment. This restriction was violated most often for memory-intensive models on datasets with many observations (see Section 4 and Appendix F.3). The tuning process is repeated independently for each tuning measure (see Section 3.4).

¹https://github.com/slds-lmu/paper_2023_survival_benchmark

Prediction types: In general, there are four prediction types in survival analysis [91]: A linear predictor `lp`, continuous ranking `crank` (e.g. a relative risk), a distribution `distr` (e.g. the survival probability), and predicted survival times `response`. The `response` time is very uncommon and only directly provided by the Survival SVM at the time of writing. Here we focus on evaluating distribution and continuous rank predictions. However, the prediction types provided by individual methods (and implementation) can vary, which is why *mlr3proba* [89] compositors are used to derive missing prediction types needed for evaluation, if necessary. Where models only predict a probabilistic prediction, `crank` is calculated as the expected mortality derived from the `distr` prediction [48, 86]. When models predict only `crank` or `lp` (i.e., RRT and SSVM) then the prediction distribution is composed with a Kaplan-Meier baseline and a PH or AFT model form, which are chosen as these are common model assumptions (which, however, doesn't seem to yield good predictions based on our results; cf. Appendix F). The XGBoost model with the Cox optimization criterion predicts `lp`, the `distr` prediction is composed using the Breslow estimator [69] similarly to previous benchmarks [49].

Pre-processing is applied only if either technically required to run a model or in line with standard recommendations for that model class. This includes standardization of covariates to unit variance and zero mean and/or treatment encoding of categorical features. We created pipelines for all learners (where required) with *mlr3pipelines* [13], to combine the respective pre-processing operation with the learning algorithm, and to properly embed the pre-processing into CV. Appendix D lists the model-specific pre-processing performed. In addition to these model-specific pre-processing steps, we collapse levels of categorical variables with frequencies below 5% as part of the model pipeline, ensuring that high-cardinality categorical features are handled consistently. As we only applied basic pre-processing, no additional hyperparameters were added to the tuning search space.

3.4 Performance Evaluation

We assess performance using two primary measures alongside seven additional measures. For cases where individual model predictions were not possible during the inner- or outer resampling procedure due to any kind of error (e.g. numerical issues), the prediction of a KM estimator was used as a fallback. This ensures a statistically sound evaluation, and is a good compromise between either assigning failed models a performance of 0 (which overpenalizes) or simply disregarding failed folds during evaluation (which is overly optimistic) [31].

Measures chosen for this benchmark are summarized in Table 1. Of these measures, only three are used to provide primary results: Harrell's C [42] for discrimination as well as the Right-Censored Log Loss (RCLL) [6, 80] and Integrated Survival Brier Score (ISBS) [38] for overall predictive ability. The benchmark procedure is run twice, each tuning either for discrimination (Harrell's C) or overall predictive ability (RCLL) and evaluated according to similar measures, while we use ISBS in addition to evaluate models tuned on either measure for comparison purposes. We also include the strictly proper scoring rules Re-Weighted Integrated Survival Log-Likelihood (RISLL) and Re-Weighted Negative Log Loss (RNLL) [90]. Furthermore, we explore the calibration measures D-Calibration [41] and Houwelingen's [47] α in Section G.3.

Statistical Analysis is conducted following Demšar [27], initially performing global Friedman rank sum tests for all measures, where the "groups" are the models and the "blocks" are the independent datasets. Significance after Bonferroni-Holm adjustment determines whether post-hoc tests are conducted. Post-hoc (multiple-testing corrected) Bonferroni-Dunn tests are conducted and presented as critical difference diagrams, using CPH as the reference model for comparison.

3.5 Datasets

To obtain a collection of suitable datasets, we ran a search across the CRAN Task View "Survival Analysis"², Python's *pycox* library and related literature (see Section 2) and existing collection of survival datasets [30], yielding over 120 datasets. After applying the dataset inclusion criteria (see Section 3.1) and removing duplicates and derivations of other datasets, a total of 32 datasets remained. Minor changes are made to variable names, recoding of factor levels, and deletion of non-informative or "illegal" covariates (e.g. ID numbers). Observations are deleted if their event time is equal to zero. Since this paper is not concerned with a model's ability to handle or impute missing covariate

²<https://cran.r-project.org/view=Survival>

Table 1: Considered performance measures. Column 1: The measure used for evaluation. Harrell’s C, RCLL, and ISBS are used for primary analysis with remaining results in Appendix F. Column 2: The corresponding measure used for tuning, i.e., models are evaluated with (1) if they were tuned on (2). Column 3: The type of evaluation measure. Column 4: Evaluated prediction type in *mlr3proba*.

Evaluation Measure	Tuning Measure	Type	Evaluates
Harrell’s C	Harrell’s C	Discrimination	crank
Uno’s C	Harrell’s C	Discrimination	crank
Integrated Survival Brier Score (ISBS)	Harrell’s C, RCLL	Scoring Rule	distr
Right-Censored Log-Likelihood (RCLL)	RCLL	Scoring Rule	distr
Re-Weighted Integrated Survival Log-Likelihood (RISLL)	RCLL	Scoring Rule	distr
Re-Weighted Negative Log Loss (RNLL)	RCLL	Scoring Rule	distr
D-Calibration	RCLL	Calibration	distr
van Houwelingen’s α	RCLL	Calibration	distr

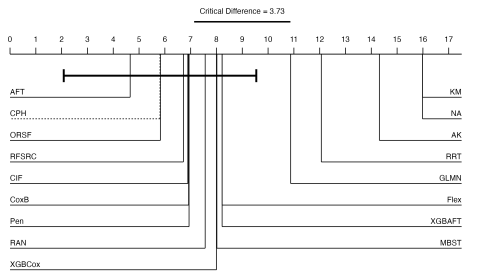
data, missing observations are removed, which occurred rarely and in small proportions. While this could introduce bias if missingness is dependent on the features or outcome, the goal of this study is model comparison, so this should not affect the conclusions. Lastly, in cases where datasets had a large number of unique time points, the time variable was coarsened via appropriate rounding, which greatly reduced computational coast for some methods including RAN, AK, MBST and CIF. For full details, see the pre-processing code contained in the GitHub repository (see Section 3.1). Summaries of the datasets in terms of the number of observations and covariates after modification and censoring proportions, along with citations for the respective sources, can be found in Appendix B.

4 Results

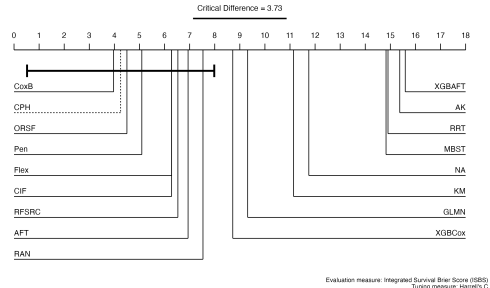
Global Friedman tests were significant for all measures, indicating the presence of significant differences between models and allowing for post-hoc analysis. Results for SSVMs are omitted due to persistent numerical and technical issues that prevented proper analysis of the algorithm’s performance. The number of times models failed to compute results due to either time or memory constraints and required the score imputation using KM (Section 3.3) is tabulated in Appendix F.3 We present critical difference (CD) plots for the baseline-comparison to CPH for both discrimination and overall performance (see Section 3.4, Demšar [27]). The top line represents a model’s average performance rank across all datasets in the benchmark, where lower ranking scores imply better performance regardless of the evaluation measure applied. Thick horizontal lines around the CPH model rank indicate the symmetric critical differences, meaning that other models within this range do not significantly differ in rank from the reference model.

Discrimination CD plots for discrimination, tuned and evaluated on Harrell’s C (Figure 1a), indicate that all models outperform the baselines learners (KM, NA, AK) except for RRT and GLMN. AFT and ORSF are the top-performing models but fail to significantly outperform the CPH baseline. The nine remaining models mostly belong to the classes of RSFs and GBMs in addition to Pen, which all achieve average ranks between 6.5 and 8.5, indicating similar discrimination performance.

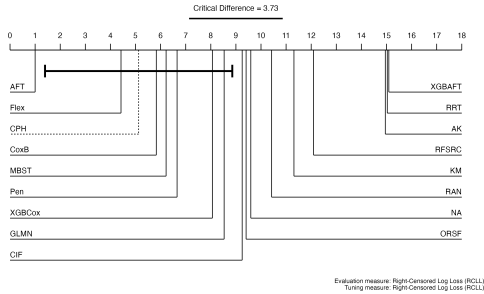
Overall Performance CD plots for overall performance evaluated on RCLL (Figure 1c) indicate that only the tuned AFT significantly outperforms the CPH model. Flex ranks better than CPH but not significantly. CPH outperforms CoxB, MBST, Pen, XGBCox, and GLMN, which do not score significantly different ranks. CIF, ORSF, RAN and RSFRC rank worse here compared to the discrimination results and are significantly outperformed by CPH. Notably, NA and KM rank better on overall performance than RFSRC, AK, RRT and XGBAFT, while RAN ranks between NA and KM. Models tuned on RCLL and evaluated on ISBS (Figure 1d) also show CPH to be among the top performing models, in this case not significantly outperformed by any other model. CoxB achieves the best average rank, with ORSF closely behind on par with CPH. AFT here performs slightly worse yet still in the top half of all models. Pen, Flex and CIF also do not significantly differ from CPH, whereas the majority of models are outperformed by CPH, from GLMN and MBST to XGBCox, RAN, KM, RFSRC and NA. XGBAFT, AK and RRT are unambiguously placed at the bottom ranks. For models tuned with Harrell’s C and evaluated on ISBS (Figure 1b), we observe a



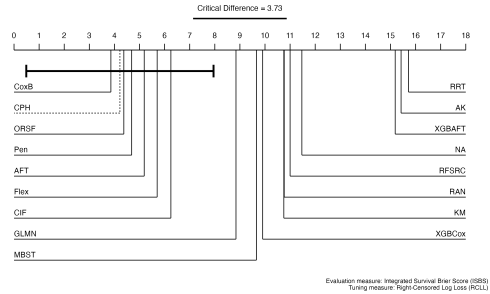
(a) Models tuned and evaluated with Harrell's C



(b) Models tuned with Harrell's C, evaluated with ISBS



(c) Models tuned and evaluated with RCLL



(d) Models tuned on RCLL, evaluated with ISBS

Figure 1: Critical difference plot comparing models with the CPH reference tuned on Harrell's C (a,b) and RCLL (c,d) and evaluated on Harrell's C (a), RCLL (c) and ISBS (b,d). Superior models (lower ranking scores) are on the left with decreasing performance (higher rank) moving right. Models connected by thick horizontal lines are not significantly different from the baseline when adjusting for multiple comparisons.

Right-Censored Log-Loss (RCLL)

Boxplot of aggregated scores across all tasks transformations

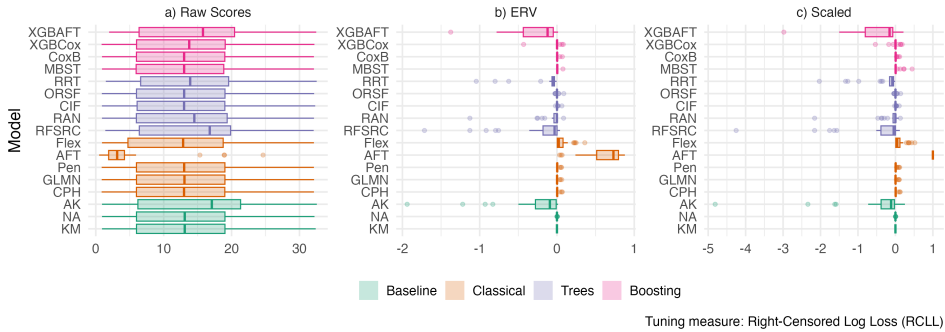


Figure 2: Boxplots of aggregated scores across all datasets for models tuned and evaluated with RCLL showing unmodified RCLL scores (a), Explained Residual Variation (ERV) scores (b), and scores scaled such that 0 is equivalent to KM and 1 is achieved by the best model for each dataset and measure.

similar ranking with CoxB achieving the best ranking score, followed by CPH, ORSF, Pen, Flex and CIF. RSFRC, AFT, and RAN complete the range of models not performing significantly different than CPH, whereas GLMN and the remaining GBMs, RSFs and baseline models perform significantly worse than CPH.

We additionally present boxplots both for individual scores per dataset and aggregated scores. We offer three versions of these aggregated scores to support evaluation and analysis, illustrated by Figure 2: a) Raw scores as calculated by the corresponding evaluation measure; b) "Explained Residual Variation" (ERV) [60] scores (similar to the "Index of Prediction Accuracy" [54]) where negative

values imply worse than KM performance, 0 is equivalent to KM, higher is better, and 1 denotes a perfect model; c) Scaled scores, whose interpretation is the same as the ERV ones, with the difference that 1 is achieved by the best model for a given task and measure [24]. Further results for all measures can be found in Appendix F. Additional tables and visualizations are available on our results website which is linked in our GitHub repository. The repository also provides downloads for all raw results and tuning archives.

5 Discussion

Discrimination The critical difference plots indicate that AFT, ORSF, and CPH were among the best performing methods, notably including two classical methods. The best-performing ML methods were ORSF ahead of CPH and RFSRC, CIF, and CoxB behind it, highlighting RSFs and GBMs good performance while not being able to significantly outperform the CPH reference. XGBoost with either Cox or AFT objective appears at the lower end of the models ranking on par with CPH, while ranking lower than the likelihood-boosting method CoxB and all RSFs. Conversely, CoxB did not require any explicit tuning outside of its internal optimization method, making it both computationally more efficient in this regard while achieving comparable or better discrimination performance. Penalized Cox (GLMN) and decision trees (RRT) are the only non-baseline models clearly outperformed by CPH, which in either case may be expected given that RRT is known to be inferior to boosted or bagged trees for prediction purposes, and GLMN representing a penalized version of CPH that may not be well suited for the low-dimensional tasks included in this benchmark. AK performs as expected, as it represents a slight improvement over KM and NA as AK takes the covariate-dependent censoring distribution into account. KM and NA are only included as a reference due to their Harrell’s C scores being constant at 0.5 by design.

Overall Performance Overall performance is judged primarily by the strictly proper scoring rule RCLL and additionally by the ISBS scoring rule [90] and should capture both a model’s discrimination and calibration properties. Notably both in RCLL and ISBS, the classical statistical methods (AFT, Flex, CPH) perform very well, with CoxB being the only ML approach consistently scoring among the top ranks for most evaluation metrics. CPH was only significantly surpassed by AFT when evaluated on RCLL, where AFT was superior by a wide margin. We note that our use of the parametric AFT model can be considered unconventional, as we tune the functional form (i.e. whether to use a Weibull, log-normal, or log-logistic distribution) within the three inner resampling folds, which lead to tuning results averaged across the five outer resampling folds where different functional forms were chosen by the tuning procedure. This also affects XGBAFT analogously. In real-world scenarios, a specific functional form is typically chosen in advance, leading us to believe that the performance of the parametric models may not be fully representative. We also note that many ML methods that perform on par with CPH on discrimination measures noticeably rank lower when considering overall performance measures. This implies that calibration is an important factor that may be neglected by some models, but where AFT is particularly strong. The difference in rankings between RCLL and ISBS stands out, as ORSF ranks very well on RCLL yet is significantly outranked by CPH on ISBS. Interestingly, the difference in rankings based on ISBS is very similar for models tuned on either Harrell’s C or RCLL, with CoxB, CPH and ORSF leading the ranking. While RCLL and ISBS are both scoring rules, they evaluate different properties of the predicted targets and therefore produce different results. Generally speaking, multiple measures should be considered for performance evaluation in any case as they may highlight aspects of performance relevant in different contexts.

Calibration In the case of D-Calibration, a model is considered well-calibrated if the underlying Pearson’s χ^2 -test results in $p > 0.05$. For van Houwelingen’s α , calibration is indicated by values close to 1. All models aside from notable outliers AK, MBST, XGBCox and XGBAFT appear to be reasonably well calibrated based on both D-Calibration and van Houwelingen’s α across most datasets (Section G.3). This is consistent with the comparatively poor performance displayed in the results evaluating on RCLL and ISBS, where RRT, AK and XGBAFT were among the lowest-scoring models. However, both measures should be considered experimental, and more research on calibration-specific measures is required for a more conclusive evaluation (e.g. [5]) For discrimination and overall performance, results on individual datasets are not discussed here but presented in Appendix G as the overall trend is similar to the aggregated results.

5.1 Limitations

We apply random search with a comparatively low evaluation budget as a tuning strategy for our benchmark, which is not necessarily the best option for some of our slightly higher-dimensional search spaces. In particular for larger search spaces, Bayesian optimization [34] could be more appropriate. However, we note that we greatly exceed the amount of tuning performed by many of the existing benchmarks in the literature (see Section 2). We distinguish between the ‘Cox’ and ‘AFT’ forms for XGBoost and treat them as separate models for the purposes of our benchmark, as the difference in objective represents very different model assumptions. Analogously, it could be argued that MBST should also be split in a similar fashion based on the family parameter. Furthermore, our use of the AFT model prioritizes prediction over interpretation and could be simplified by splitting the tuned model into one model for each functional form (“Weibull”, “log-normal”, “log-logistic”), which would yield more interpretable results. Finally, since our focus lies on generalizability to low-dimensional, right-censored settings, our results will not generalize to more complex settings. However, extension to left-censoring, competing risks or other more complex endpoints first necessitates more comprehensive support by models and their implementations as well as adaptation of available evaluation measures to these scenarios. The number of datasets included could be extended, but still exceeds that of the vast majority of previous benchmarks.

5.2 Conclusions and Future Work

Our results demonstrate that classical statistical methods, such as CPH and AFT, can significantly outperform complex ML algorithms in predictive survival tasks. While it is possible to fine-tune and achieve better predictive accuracy using ML methods in individual cases (see Appendix G), our results indicate that across a range of low-dimensional tasks, this is not the case in aggregate. We therefore recommend practitioners to start with these conceptually and computationally simpler methods, and evaluate whether the additional computational cost and loss of interpretability is appropriate for their needs. To improve upon our results, future work might employ alternative tuning strategies such as Bayesian Optimization or different model configurations with adjusted hyperparameter spaces. Additionally, an expansion of this benchmark with a wider range of settings would be beneficial, such that results can be generalized further pending the corresponding software support.

Acknowledgments and Disclosure of Funding

We are grateful for the package maintainers' time supporting this effort regarding the model configurations used in this benchmark.

We gratefully acknowledge that a pilot version of this study was performed by Sonabend [91], laying the ground work for the experiments conducted in this study.

This work has been carried out by making use of Wyoming's Advanced Research Computing Center, on its Beartooth Compute Environment (<https://doi.org/10.15786/M2FY47>). We gratefully acknowledge the computational and data resources provided by Wyoming's Advanced Research Computing Center (<https://www.uwyo.edu/arcc/>).

JZ received funding from the European Union's Horizon 2020 research and innovation program under grant agreement No 101016851, project PANCAIM.

BB is supported by the German Research Foundation (DFG), Grant Number: 460135501.

MNW is supported by the German Research Foundation (DFG), Grant Numbers: 437611051, 459360854.

References

- [1] Odd Aalen. "Nonparametric Inference for a Family of Counting Processes". In: *The Annals of Statistics* 6.4 (1978), pp. 701–726.
- [2] Michael G Akritas. "Nearest Neighbor Estimation of a Bivariate Distribution Under Random Censoring". en. In: *Ann. Statist.* 22.3 (1994), pp. 1299–1327. ISSN: 0090-5364. DOI: 10.1214/aos/1176325630. URL: <https://projecteuclid.org:443/euclid.aos/1176325630>.
- [3] Alina M Allen et al. "Nonalcoholic fatty liver disease incidence and impact on metabolic burden and death: A 20 year-community study." eng. In: *Hepatology (Baltimore, Md.)* 67.5 (Apr. 2018), pp. 1726–1736. ISSN: 1527-3350 (Electronic). DOI: 10.1002/hep.29546.
- [4] Per Kragh Andersen et al. "Introduction". In: *Statistical Models Based on Counting Processes*. Springer Series in Statistics. New York, NY: Springer US, 1993, pp. 1–44. DOI: <https://doi.org/10.1007/978-1-4612-4348-9>.
- [5] Peter C. Austin, Frank E. Harrell Jr, and David van Klaveren. "Graphical Calibration Curves and the Integrated Calibration Index (ICI) for Survival Models". In: *Statistics in Medicine* 39.21 (2020), pp. 2714–2742. ISSN: 1097-0258. DOI: 10.1002/sim.8570.
- [6] Anand Avati et al. "Countdown Regression: Sharp and Calibrated Survival Predictions". In: *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*. Uncertainty in Artificial Intelligence. PMLR, Aug. 6, 2020, pp. 145–155. URL: <https://proceedings.mlr.press/v115/avati20a.html> (visited on 05/21/2024).
- [7] Avinash Barnwal, Hyunsu Cho, and Toby Hocking. "Survival Regression with Accelerated Failure Time Model in XGBoost". In: *Journal of Computational and Graphical Statistics* 31.4 (Oct. 2, 2022), pp. 1292–1302. ISSN: 1061-8600. DOI: 10.1080/10618600.2022.2067548.
- [8] Beartooth Computing Environment, x86_64 cluster. *Advanced Research Computing Center*. University of Wyoming, Laramie, WY. 2024. DOI: <https://doi.org/10.15786/M2FY47>.
- [9] Marc Becker. *mlr3tuningspaces: Search Spaces for 'mlr3'*. GitHub, 2024. URL: <https://github.com/mlr-org/mlr3extralearners>.
- [10] Andreas Bender and Fabian Scheipl. "pamtools: Piece-wise exponential Additive Mixed Modeling tools". In: *arXiv:1806.01042 [stat]* (2018). URL: <http://arxiv.org/abs/1806.01042>.
- [11] Andreas Bender et al. "Penalized estimation of complex, non-linear exposure-lag-response associations". In: *Biostatistics* 20.2 (Feb. 2018), pp. 315–331. ISSN: 1465-4644. DOI: 10.1093/biostatistics/kxy003. URL: <https://doi.org/10.1093/biostatistics/kxy003>.
- [12] Harald Binder and Martin Schumacher. "Allowing for Mandatory Covariates in Boosting Estimation of Sparse High-Dimensional Survival Models". In: *BMC Bioinformatics* 9.1 (Jan. 10, 2008), p. 14. ISSN: 1471-2105. DOI: 10.1186/1471-2105-9-14.

- [13] Martin Binder et al. “mlr3pipelines - Flexible Machine Learning Pipelines in R”. In: *Journal of Machine Learning Research* 22.184 (2021), pp. 1–7. URL: <http://jmlr.org/papers/v22/21-0281.html>.
- [14] Bernd Bischl et al. “Hyperparameter Optimization: Foundations, Algorithms, Best Practices, and Open Challenges”. In: *WIREs Data Mining and Knowledge Discovery* 13.2 (2023), e1484. ISSN: 1942-4795. DOI: 10.1002/widm.1484.
- [15] Paul Blanche, Jean-François Dartigues, and H el ene Jacquemin-Gadda. “Estimating and Comparing Time-Dependent Areas under Receiver Operating Characteristic Curves for Censored Event Times with Competing Risks”. In: *Statistics in Medicine* 32.30 (2013), pp. 5381–5397. ISSN: 1097-0258. DOI: 10.1002/sim.5958.
- [16] Anne-Laure Boulesteix, Sabine Lauer, and Manuel J. A. Eugster. “A Plea for Neutral Comparison Studies in Computational Sciences”. In: *PLOS ONE* 8.4 (Apr. 24, 2013), e61562. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0061562.
- [17] Leo Breiman, ed. *Classification and Regression Trees*. 1. CRC Press repr. Boca Raton, Fla.: Chapman & Hall/CRC, 1998. 358 pp. ISBN: 978-0-412-04841-8.
- [18] N. E. Breslow and N. Chatterjee. “Design and Analysis of Two-Phase Studies with Binary Outcome Applied to Wilms Tumour Prognosis”. In: *Journal of the Royal Statistical Society Series C: Applied Statistics* 48.4 (Dec. 1, 1999), pp. 457–468. ISSN: 0035-9254. DOI: 10.1111/1467-9876.00165.
- [19] G oran Brostr om. *eha: Event History Analysis*. R package version 2.9.0. 2021. URL: <http://ehar.se/r/eha/>.
- [20] Peter B uhlmann and Bin Yu. “Boosting With the L2 Loss: Regression and Classification”. In: *Journal of the American Statistical Association* 98.462 (June 1, 2003), pp. 324–339. ISSN: 0162-1459. DOI: 10.1198/016214503000125.
- [21] Chao Cai et al. *smcure: Fit Semiparametric Mixture Cure Models*. R package version 2.0. 2012. URL: <https://CRAN.R-project.org/package=smcure>.
- [22] Bradley P. Carlin and Thomas A. Louis. “Supplemental Materials to Bayesian Methods for Data Analysis, 3rd Edition”. In: (Oct. 2, 2018). DOI: 10.13020/D6N10N.
- [23] Daniel P. Carpenter. “Groups, the Media, Agency Waiting Costs, and FDA Drug Approval”. In: *American Journal of Political Science* 46.3 (2002), pp. 490–505. ISSN: 00925853, 15405907. URL: <http://www.jstor.org/stable/3088394>.
- [24] Rich Caruana and Alexandru Niculescu-Mizil. “An Empirical Comparison of Supervised Learning Algorithms”. In: *Proceedings of the 23rd International Conference on Machine Learning - ICML ’06*. The 23rd International Conference. Pittsburgh, Pennsylvania: ACM Press, 2006, pp. 161–168. ISBN: 978-1-59593-383-6. DOI: 10.1145/1143844.1143865.
- [25] Tianqi Chen and Carlos Guestrin. “XGBoost: A Scalable Tree Boosting System”. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD ’16. New York, NY, USA: Association for Computing Machinery, Aug. 13, 2016, pp. 785–794. ISBN: 978-1-4503-4232-2. DOI: 10.1145/2939672.2939785.
- [26] D. R. Cox. “Partial Likelihood”. In: *Biometrika* 62.2 (Aug. 1, 1975), pp. 269–276. ISSN: 0006-3444. DOI: 10.1093/biomet/62.2.269.
- [27] Janez Dem sar. “Statistical comparisons of classifiers over multiple data sets”. In: *Journal of Machine learning research* 7.1 (2006), pp. 1–30.
- [28] Lore Dirick, Gerda Claeskens, and Bart Baesens. “Time to default in credit scoring using survival analysis: A benchmark study”. In: *Journal of the Operational Research Society* 68.6 (2017), pp. 652–665. ISSN: 14769360. DOI: 10.1057/s41274-016-0128-9.
- [29] Angela Dispenzieri et al. “Use of nonclonal serum immunoglobulin free light chains to predict overall survival in the general population”. eng. In: *Mayo Clinic proceedings* 87.6 (June 2012), pp. 517–523. ISSN: 1942-5546. DOI: 10.1016/j.mayocp.2012.03.009. URL: <https://www.ncbi.nlm.nih.gov/pubmed/22677072><https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3538473/>.
- [30] Erik Drysdale. “SurvSet: An Open-Source Time-to-Event Dataset Repository”. Mar. 6, 2022. arXiv: 2203.03094 [cs, stat]. URL: <http://arxiv.org/abs/2203.03094> (visited on 03/12/2022).
- [31] Sebastian Fischer, Michel Lang, and Marc Becker. “Large-Scale Benchmarking”. In: *Applied Machine Learning Using mlr3 in R*. Ed. by Bernd Bischl et al. CRC Press, 2024. URL: https://mlr3book.mlr-org.com/large-scale_benchmarking.html.

- [32] Y. Foucher et al. *RISCA: Causal Inference and Prediction in Cohort-Based Analyses*. R package version 1.0.4. 2023. URL: <https://CRAN.R-project.org/package=RISCA>.
- [33] Jerome Friedman. “Stochastic Gradient Boosting”. In: *Computational Statistics & Data Analysis* 38 (Mar. 1999), pp. 367–378. DOI: 10.1016/S0167-9473(01)00065-2.
- [34] Roman Garnett. *Bayesian Optimization*. Cambridge University Press, 2023.
- [35] Ekavi N Georgousopoulou et al. “Comparisons between Survival Models in Predicting Cardiovascular Disease Events : Application in the ATTICA Study (2002-2012).” In: *Journal of Statistics Applications & Probability* 4.2 (2015), pp. 203–210.
- [36] Jelle J. Goeman. “L1 Penalized Estimation in the Cox Proportional Hazards Model”. In: *Biometrical Journal* 52.1 (2010), pp. 70–84. ISSN: 1521-4036. DOI: 10.1002/bimj.200900028.
- [37] Shahrbanoo Goli et al. “Performance Evaluation of Support Vector Regression Models for Survival Analysis: A Simulation Study”. In: *International Journal of Advanced Computer Science and Applications* 7.6 (2016). ISSN: 21565570, 2158107X. DOI: 10.14569/IJACSA.2016.070650.
- [38] E. Graf et al. “Assessment and Comparison of Prognostic Classification Schemes for Survival Data”. In: *Statistics in Medicine* 18.17-18 (Sept. 15, 1999), pp. 2529–2545. ISSN: 0277-6715. DOI: 10.1002/(sici)1097-0258(19990915/30)18:17/18<2529::aid-sim274>3.0.co;2-5. pmid: 10474158.
- [39] Il Do Ha et al. *frailtyHL: Frailty Models via Hierarchical Likelihood*. R package version 2.3. 2019. URL: <https://CRAN.R-project.org/package=frailtyHL>.
- [40] Danial Habibi et al. “Comparison of Survival Models for Analyzing Prognostic Factors in Gastric Cancer Patients”. In: *Asian Pacific journal of cancer prevention : APJCP* 19.3 (2018), pp. 749–753. ISSN: 2476-762X. DOI: 10.22034/APJCP.2018.19.3.749. URL: <http://www.ncbi.nlm.nih.gov/pubmed/29582630> <http://www.ncbi.nlm.nih.gov/pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC5980851>.
- [41] Humza Haider et al. “Effective Ways to Build and Evaluate Individual Survival Distributions”. In: *Journal of Machine Learning Research* 21.85 (2020), pp. 1–63. ISSN: 1533-7928. URL: <http://jmlr.org/papers/v21/18-772.html> (visited on 05/21/2024).
- [42] Frank E. Harrell, Robert M. Califf, and David B. Pryor. “Evaluating the yield of medical tests”. In: *JAMA* 247.18 (May 1982), pp. 2543–2546. ISSN: 0098-7484. URL: <http://dx.doi.org/10.1001/jama.1982.03320430047030>.
- [43] Moritz Herrmann et al. “Large-Scale Benchmark Study of Survival Prediction Methods Using Multi-Omics Data”. In: *Briefings in Bioinformatics* 22.3 (May 1, 2021), bbaa167. ISSN: 1477-4054. DOI: 10.1093/bib/bbaa167.
- [44] David W Hosmer, Stanley Lemeshow, and Susanne May. *Applied survival analysis regression modeling of time-to-event data*. eng. 2nd ed. Wiley series in probability and statistics. Hoboken, N.J.: Wiley-Interscience, 2008.
- [45] David W Hosmer Jr, Stanley Lemeshow, and Susanne May. *Applied survival analysis: regression modeling of time-to-event data*. Vol. 618. John Wiley & Sons, 2011.
- [46] Torsten Hothorn, Kurt Hornik, and Achim Zeileis. “Unbiased Recursive Partitioning: A Conditional Inference Framework”. In: *Journal of Computational and Graphical Statistics* 15.3 (Sept. 1, 2006), pp. 651–674. ISSN: 1061-8600. DOI: 10.1198/106186006X133933.
- [47] Hans C. van Houwelingen. “Validation, Calibration, Revision and Combination of Prognostic Survival Models”. In: *Statistics in Medicine* 19.24 (2000), pp. 3401–3415. ISSN: 1097-0258. DOI: 10.1002/1097-0258(20001230)19:24<3401::AID-SIM554>3.0.CO;2-2.
- [48] By Hemant Ishwaran et al. “Random survival forests”. In: *The Annals of Statistics* 2.3 (2008), pp. 841–860. DOI: 10.1214/08-AOAS169. arXiv: arXiv:0811.1645v1.
- [49] Byron C. Jaeger et al. “Accelerated and Interpretable Oblique Random Survival Forests”. In: *Journal of Computational and Graphical Statistics* 33.1 (Jan. 2, 2024), pp. 192–207. ISSN: 1061-8600. DOI: 10.1080/10618600.2023.2231048.
- [50] Henrik Stig Jørgensen et al. “Acute Stroke With Atrial Fibrillation”. In: *Stroke* 27.10 (Oct. 1, 1996), pp. 1765–1769. DOI: 10.1161/01.STR.27.10.1765.
- [51] John D. Kalbfleisch and Ross L. Prentice. *The Statistical Analysis of Failure Time Data*. John Wiley & Sons, Jan. 25, 2011. 464 pp. ISBN: 978-1-118-03123-0.

- [52] E. L. Kaplan and Paul Meier. “Nonparametric Estimation from Incomplete Observations”. In: *Journal of the American Statistical Association* 53.282 (1958), pp. 457–481. ISSN: 01621459. DOI: 10.2307/2281868.
- [53] Michael W. Kattan. “Comparison of Cox Regression With Other Methods for Determining Prediction Models and Nomograms”. In: *The Journal of Urology*. Part 2 of 2 170 (6, Supplement Dec. 1, 2003), S6–S10. ISSN: 0022-5347. DOI: 10.1097/01.ju.0000094764.56269.2d.
- [54] Michael W. Kattan and Thomas A. Gerds. “The Index of Prediction Accuracy: An Intuitive Measure Useful for Evaluating Risk Prediction Models”. In: *Diagnostic and Prognostic Research* 2.1 (May 4, 2018), p. 7. ISSN: 2397-7523. DOI: 10.1186/s41512-018-0029-2.
- [55] Jared Katzman. *DeepSurv*. 2016. URL: <https://pypi.org/project/deepsurv/>.
- [56] Jared L Katzman et al. “DeepSurv: personalized treatment recommender system using a Cox proportional hazards deep neural network”. In: *BMC Medical Research Methodology* 18.1 (2018), p. 24. ISSN: 1471-2288. DOI: 10.1186/s12874-018-0482-1. URL: <https://doi.org/10.1186/s12874-018-0482-1>.
- [57] John M Kirkwood et al. “Interferon alfa-2b adjuvant therapy of high-risk resected cutaneous melanoma: the Eastern Cooperative Oncology Group Trial EST 1684.” In: *Journal of clinical oncology* 14.1 (1996), pp. 7–17.
- [58] John P Klein and Melvin L Moeschberger. *Survival analysis: techniques for censored and truncated data*. 2nd ed. Springer Science & Business Media, 2003. ISBN: 0387216456.
- [59] Roger Koenker. *quantreg: Quantile Regression*. R package version 5.86. 2021. URL: <https://www.r-project.org>.
- [60] Edward L Korn and Richard Simon. “Explained Residual Variation, Explained Risk, and Goodness of Fit”. In: *The American Statistician* 45.3 (Dec. 1991), pp. 201–206. ISSN: 00031305. DOI: 10.2307/2684290. URL: <http://www.jstor.org/stable/2684290>.
- [61] Håvard Kvamme. *pycox*. 2018. URL: <https://pypi.org/project/pycox/>.
- [62] R A Kyle. ““Benign” monoclonal gammopathy—after 20 to 35 years of follow-up.” eng. In: *Mayo Clinic proceedings* 68.1 (Jan. 1993), pp. 26–36. ISSN: 0025-6196 (Print). DOI: 10.1016/s0025-6196(12)60015-9.
- [63] Pierre L’Ecuyer. “Good Parameters and Implementations for Combined Multiple Recursive Random Number Generators”. In: *Operations Research* 47.1 (Feb. 1999), pp. 159–164. ISSN: 0030-364X. DOI: 10.1287/opre.47.1.159. URL: <https://pubsonline.informs.org/doi/abs/10.1287/opre.47.1.159>.
- [64] Michel Lang, Bernd Bischl, and Dirk Surmann. “Batchtools: Tools for R to Work on Batch Systems”. In: *Journal of Open Source Software* 2.10 (Feb. 22, 2017), p. 135. ISSN: 2475-9066. DOI: 10.21105/joss.00135.
- [65] Michel Lang et al. “mlr3: A modern object-oriented machine learning framework in R”. In: *Journal of Open Source Software* 4.44 (2019), p. 1903. DOI: 10.21105/joss.01903. URL: <https://joss.theoj.org/papers/10.21105/joss.01903>.
- [66] Michel Lang et al. *mlr3tuning: Tuning for ‘mlr3’*. 2023. URL: <https://cran.r-project.org/package=mlr3tuning>.
- [67] Seungyeoun Lee and Heeju Lim. “Review of statistical methods for survival analysis using genomic data”. eng. In: *Genomics & informatics* 17.4 (Dec. 2019), e41–e41. ISSN: 1598-866X. DOI: 10.5808/GI.2019.17.4.e41. URL: <https://pubmed.ncbi.nlm.nih.gov/31896241/> <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6944043/>.
- [68] Lee, Changhee et al. *DeepHit*. 2019. URL: <https://github.com/ch18856/DeepHit>.
- [69] D. Y. Lin. “On the Breslow estimator”. In: *Lifetime Data Analysis* 13.4 (2007), pp. 471–480. ISSN: 13807870. DOI: 10.1007/s10985-007-9048-y.
- [70] C L Loprinzi et al. “Prospective evaluation of prognostic variables from patient-completed questionnaires. North Central Cancer Treatment Group.” eng. In: *Journal of clinical oncology : official journal of the American Society of Clinical Oncology* 12.3 (Mar. 1994), pp. 601–607. ISSN: 0732-183X (Print). DOI: 10.1200/JCO.1994.12.3.601.
- [71] James T. Luxhoj and Huan-Jyh Shyur. “Comparison of Proportional Hazards Models and Neural Networks for Reliability Estimation”. In: *Journal of Intelligent Manufacturing* 8.3 (May 1, 1997), pp. 227–234. ISSN: 1572-8145. DOI: 10.1023/A:1018525308809.

- [72] M. Pohar and J. Stare. “Relative survival analysis in R”. In: *Computer methods and programs in biomedicine* 81 (3 2006), pp. 272–278. DOI: 10.1016/j.cmpb.2006.01.004.
- [73] Ulla B Mogensen, Hemant Ishwaran, and Thomas A Gerds. *Evaluating Random Forests for Survival Analysis using Prediction Error Curves*. 2014.
- [74] Bijan Moghimi-Dehkordi et al. “Statistical Comparison of Survival Models for Analysis of Cancer Data”. In: *Asian Pacific journal of cancer prevention: APJCP* 9.3 (2008), pp. 417–420. ISSN: 2476-762X. PMID: 18990013.
- [75] John V. Monaco, Malka Gorfine, and Li Hsu. “General Semiparametric Shared Frailty Model: Estimation and Simulation with frailtySurv”. In: *Journal of Statistical Software* 86.4 (2018), pp. 1–42. DOI: 10.18637/jss.v086.i04.
- [76] Lucila Ohno-Machado. “A Comparison of Cox Proportional Hazards and Artificial Neural Network Models for Medical Prognosis”. In: *Computers in Biology and Medicine* 27.1 (Jan. 1, 1997), pp. 55–65. ISSN: 0010-4825. DOI: 10.1016/S0010-4825(96)00036-4.
- [77] Lucila Ohno-Machado. “Modeling medical prognosis: survival analysis techniques.” eng. In: *Journal of biomedical informatics* 34.6 (Dec. 2001), pp. 428–439. ISSN: 1532-0464 (Print). DOI: 10.1006/jbin.2002.1038.
- [78] Katie Patel, Richard Kay, and Lucy Rowell. “Comparing Proportional Hazards and Accelerated Failure Time Models: An Application in Influenza”. In: *Pharmaceutical Statistics* 5.3 (2006), pp. 213–224. ISSN: 1539-1612. DOI: 10.1002/pst.213.
- [79] R Core Team. *R: A Language and Environment for Statistical Computing*. Vienna, 2017.
- [80] David Rindt et al. “Survival Regression with Proper Scoring Rules and Monotonic Neural Networks”. In: *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*. International Conference on Artificial Intelligence and Statistics. PMLR, May 3, 2022, pp. 1190–1205. URL: <https://proceedings.mlr.press/v151/rindt22a.html> (visited on 05/21/2024).
- [81] Dimitris Rizopoulos. “JM: An R Package for the Joint Modelling of Longitudinal and Time-to-Event Data”. In: *Journal of Statistical Software* 35.9 (2010), pp. 1–33. URL: <http://www.jstatsoft.org/v35/i09/>.
- [82] Patrick Royston and Mahesh K.B. Parmar. “Flexible parametric proportional-hazards and proportional-odds models for censored survival data, with application to prognostic modelling and estimation of treatment effects”. In: *Statistics in Medicine* 21.15 (2002), pp. 2175–2197. ISSN: 02776715. DOI: 10.1002/sim.1203.
- [83] Matthias Schmid and Sergej Potapov. “A Comparison of Estimators to Evaluate the Discriminatory Power of Time-to-Event Models”. In: *Statistics in Medicine* 31.23 (2012), pp. 2588–2609. ISSN: 1097-0258. DOI: 10.1002/sim.5464.
- [84] Noah Simon et al. “Regularization Paths for Cox’s Proportional Hazards Model via Coordinate Descent”. In: *Journal of Statistical Software* 39 (Mar. 9, 2011), pp. 1–13. ISSN: 1548-7660. DOI: 10.18637/jss.v039.i05.
- [85] Raphael Sonabend. *survivalmodels: Models for Survival Analysis*. 2020. URL: <https://cran.r-project.org/package=survivalmodels>.
- [86] Raphael Sonabend, Andreas Bender, and Sebastian Vollmer. “Avoiding C-hacking When Evaluating Survival Distribution Predictions with Discrimination Measures”. In: *Bioinformatics* 38.17 (Sept. 2, 2022), pp. 4178–4184. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btac451.
- [87] Raphael Sonabend and Florian Pfisterer. *mlr3benchmark: Benchmarking analysis for 'mlr3'*. 2020. URL: <https://cran.r-project.org/package=mlr3benchmark>.
- [88] Raphael Sonabend and Patrick Schratz. *mlr3extralearners: Extra Learners For mlr3*. 2024. URL: <https://github.com/mlr-org/mlr3extralearners>.
- [89] Raphael Sonabend et al. “mlr3proba: An R Package for Machine Learning in Survival Analysis”. In: *Bioinformatics* (Feb. 2021). ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btab039. URL: <https://cran.r-project.org/package=mlr3proba>.
- [90] Raphael Sonabend et al. “Examining properness in the external validation of survival models with squared and logarithmic losses”. In: (2024). DOI: 10.48550/arXiv.2212.05260. arXiv: 2212.05260. URL: <https://arxiv.org/abs/2212.05260>.

- [91] Raphael Edward Benjamin Sonabend. “A Theoretical and Methodological Framework for Machine Learning in Survival Analysis: Enabling Transparent and Accessible Predictive Modelling on Right-Censored Time-to-Event Data”. Doctoral. UCL (University College London), June 28, 2021. 345 pp. URL: <https://discovery.ucl.ac.uk/id/eprint/10129352/> (visited on 06/03/2024).
- [92] Annette Spooner et al. “A comparison of machine learning methods for survival analysis of high-dimensional clinical data for dementia prediction”. In: *Scientific Reports* 10.1 (2020), p. 20410. ISSN: 2045-2322. DOI: 10.1038/s41598-020-77220-w. URL: <https://doi.org/10.1038/s41598-020-77220-w>.
- [93] Richard J. Sylvester et al. “Predicting Recurrence and Progression in Individual Patients with Stage Ta T1 Bladder Cancer Using EORTC Risk Tables: A Combined Analysis of 2596 Patients from Seven EORTC Trials”. In: *European Urology* 49.3 (2006), pp. 466–477. ISSN: 0302-2838. DOI: 10.1016/j.eururo.2005.12.031.
- [94] The Benelux C M L Study Group. “Randomized Study on Hydroxyurea Alone Versus Hydroxyurea Combined With Low-Dose Interferon- α 2b for Chronic Myeloid Leukemia”. In: *Blood* 91.8 (Apr. 1998), pp. 2713–2721. ISSN: 1528-0020. DOI: 10.1182/blood.V91.8.2713.2713_2713_2721. URL: https://doi.org/10.1182/blood.V91.8.2713.2713_2713_2721 <https://ashpublications.org/blood/article/91/8/2713/107615/Randomized-Study-on-Hydroxyurea-Alone-Versus>.
- [95] Katy Trébern-Launay et al. “Comparison of the Risk Factors Effects between Two Populations: Two Alternative Approaches Illustrated by the Analysis of First and Second Kidney Transplant Recipients”. In: *BMC Medical Research Methodology* 13 (Aug. 2013), p. 102. ISSN: 1471-2288. DOI: 10.1186/1471-2288-13-102.
- [96] Kevin Ushey and Hadley Wickham. *renv: Project Environments*. R package version 1.0.7. 2024. URL: <https://CRAN.R-project.org/package=renv>.
- [97] Vanya Van Belle et al. “Support Vector Methods for Survival Analysis: A Comparison between Ranking and Regression Approaches”. In: *Artificial Intelligence in Medicine* 53.2 (Oct. 1, 2011), pp. 107–118. ISSN: 0933-3657. DOI: 10.1016/j.artmed.2011.06.006.
- [98] J C Van Houwelingen et al. “Predictability of the survival of patients with advanced ovarian cancer.” In: *Journal of Clinical Oncology* 7.6 (1989), pp. 769–773. ISSN: 0732-183X.
- [99] Joaquin Vanschoren et al. “OpenML: Networked Science in Machine Learning”. In: *SIGKDD Explorations* 15.2 (2013), pp. 49–60. DOI: 10.1145/2641190.2641198. URL: <http://doi.acm.org/10.1145/2641190.2641198>.
- [100] William N. Venables and Brian D. Ripley. *Modern Applied Statistics with S*. 4th ed. Statistics and Computing. New York: Springer, 2002. ISBN: 978-0-387-95457-8.
- [101] Ping Wang, Yan Li, and Chandan K. Reddy. “Machine Learning for Survival Analysis: A Survey”. In: *ACM Computing Surveys* 51.6 (Feb. 27, 2019), 110:1–110:36. ISSN: 0360-0300. DOI: 10.1145/3214306.
- [102] Simon Wiegerebe et al. “Deep Learning for Survival Analysis: A Review”. In: *Artificial Intelligence Review* 57.3 (Feb. 19, 2024), p. 65. ISSN: 1573-7462. DOI: 10.1007/s10462-023-10681-3.
- [103] Paula Williamson et al. “Joint modelling of longitudinal and competing risks data.” In: *Statistics in Medicine* 27 (2008), pp. 6426–6438.
- [104] David Wissel, Daniel Rowson, and Valentina Boeva. “Systematic Comparison of Multi-Omics Survival Models Reveals a Widespread Lack of Noise Resistance”. In: *Cell Reports Methods* 3.4 (Apr. 2023), p. 100461. ISSN: 26672375. DOI: 10.1016/j.crmeth.2023.100461.
- [105] Marvin N. Wright and Andreas Ziegler. “ranger: A Fast Implementation of Random Forests for High Dimensional Data in C++ and R”. In: *Journal of Statistical Software* 77.1 (2017), pp. 1–17.
- [106] Ali Zare et al. “A Comparison between Accelerated Failure-time and Cox Proportional Hazard Models in Analyzing the Survival of Gastric Cancer Patients.” In: *Iranian journal of public health* 44.8 (2015), pp. 1095–102. ISSN: 03044556. DOI: 10.1007/s00606-006-0435-8. URL: <http://www.ncbi.nlm.nih.gov/pubmed/26587473> <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4645729>.

- [107] Yunwei Zhang et al. “SurvBenchmark: Comprehensive Benchmarking Study of Survival Analysis Methods Using Both Omics Data and Clinical Data”. en. In: *bioRxiv* (July 2021), p. 2021.07.11.451967. DOI: 10.1101/2021.07.11.451967.
- [108] Zhi Zhao et al. “Tutorial on Survival Modeling with Applications to Omics Data”. In: *Bioinformatics* 40.3 (Mar. 1, 2024), btae132. ISSN: 1367-4811. DOI: 10.1093/bioinformatics/btae132.

Appendices

A Literature Review

The following represents an extended literature review complementary to the summary provided in the main article.

Comparisons of ML and Classical Models The experiments carried out in this paper fall into this category. Only two prior experiments could be found that neutrally benchmarked more than one ML model class on low-dimensional data. Kattan (2003) [53] benchmarked tree-based models, ANNs and CPH with Harrell’s C-index across three datasets with varying censoring proportions. The models are compared for significant differences by repeating the experiments up to 50 times with different seeds thus allowing for different hyper-parameter configurations and folds in cross-validation. Boxplots across all replications indicate that no machine learning model outperformed the CPH. Zhang *et al.* (2021) [107] compare classical and ML methods, taking into account feasibility and computational efficiency for various tasks in the biomedical field. Methods are evaluated on six clinical and 16 omics datasets using 11 metrics, including time-dependent AUC, Brier score and multiple variations of the C-index. However, methods were applied with specific hyperparameter sets without tuning, thereby limiting the generalizability of their results.

Comparisons on High-Dimensional Data Herrmann *et al.* (2020) [43] performed a large-scale benchmark experiment of survival models on multi-omics high-dimensional data. Models fall into the following groups: Penalised regression, GBMs, and RSFs. Comparisons are made with Uno’s C and the Integrated Survival Brier Score (ISBS). The ISBS for all models overlapped with the Kaplan-Meier baseline though all C-indices were significantly higher than the baseline — however it is not stated how the standard errors for the confidence intervals were derived nor is it stated in the paper if multiple testing correction is applied. The authors also note that their results should be treated with caution due to the small performance differences and high variability. Spooner *et al.* (2020) [92] also compared machine learning models on high-dimensional data. In this study GBMs, RSFs, CPH and some extensions thereof were compared. Models were evaluated by Harrell’s C-index only. The results indicated that all models outperformed CPH when no additional feature selection was used but that there were no significant differences when feature selection was applied to the Cox model. There were few significant statistical differences between models. Wissel *et al.* (2023) [104] provide a systematic comparison of multi-omics cancer survival models comparing eight DL methods, RSFs and CPH. They primarily focus on the noise-resistance of these models for high-dimensional settings, finding a general lack thereof when evaluating on Antolini’s C-index and ISBS.

Comparisons of Classical Models Moghimi-dehkordi *et al.* (2008) [74] compare CPH to AFT models with various distributions. Out-of-sample measures for comparison are not provided though the AIC produced by the CPH is far higher (and therefore inferior) than those of the parametric models. Model inspection demonstrated that all models provided similar (non-significantly different) confidence intervals for hazard ratios. Georgousopoulou *et al.* (2015) [35] compared the CPH to a Weibull and Exponential AFT model. Again no out-of-sample measures were utilized, models were compared by the Cox-Snell residuals and the Bayesian Information Criterion (BIC). Similarly to Moghimi-dehkordi *et al.*, hazard ratios produced from all three models were nearly identical. The authors claim the CPH is inferior to the parametric models though only graphical comparisons are included. Zare *et al.* (2015) [106] provide another comparison of the CPH to AFT models, using Cox-Snell residuals and AIC as their measures of comparison. Similarly to the previous studies the Cox model has the highest AIC though the plotted Cox-Snell residuals are very similar. The authors acknowledge no significant differences between the model classes and instead conclude that AFT is a useful and more interpretable alternative. No significant differences were found between the different AFT parameterizations. Dirick *et al.* (2017) [28] make use of a financial setting to compare the CPH, AFT, flexible Cox models using splines to model the hazard, and mixture cure models. The authors compare the models using a time-dependent AUC, the mean squared error (MSE), and the mean absolute error (MAE). Survival times are generated from the CPH with a deterministic composition using quantiles chosen to minimize the MSE and MAE, and it is not clear if this is performed in an unbiased nested resampling manner or after predictions are made. By averaging the ranking of model

performance the authors conclude that CPH with penalized splines outperformed the other models with respect to the chosen metrics.

Habibi *et al.* (2018) [40] performed another experiment on PH and AFT models. Models were again compared exclusively by the AIC with the PH having the highest result and log-normal AFT the lowest; differences between AFT models were non-significant. Confidence intervals for hazard ratios were similar (non-significantly different) for all models.

Comparisons of a Novel Model Class Luxhoj and Shyur (1997) [71] benchmarked neural networks against CPH in the engineering field of reliability analysis. The baseline hazard of the Cox model is modeled by splines with a single knot. The models are compared using the mean squared error on a validation set of sample size 40, of these 40 there are only 9 unique failure times that are used for model testing. Insufficient information is provided to determine the architecture or training procedure of the neural networks compared. The MSE difference between the Cox and ANN was 0.003, which is highly unlikely to be a significant difference on a test set of only 40 observations with nine observed events. Ohno-Machado (1997) [76] also compared CPH to ANNs. Several Cox models were fit with automated variable selection by backwards elimination. For each model, survival curves were predicted and for a given patient they were considered dead at a particular time point if the predicted survival curve at the time is less than the “arbitrary” [76] probability of 0.5. The Cox models were compared to a single hidden layer ANN. This model made probabilistic predictions of death in four time-intervals that were within the predicted time of the Cox models. The probabilistic predictions from both models were compared with the AUC and its corresponding ROC. No significant differences in performance were found between the two models. Goli *et al.* (2016) [37] provide a comprehensive comparison of support vector machine models with CPH as a reference class (Kaplan-Meier is not included). Models are compared against the C-index and log-rank test, though it is unstated which C-index is utilized. No model outperformed CPH with respect to the chosen C-index. Jaeger *et al.* (2019) [49] compare a multiple variations of a novel implementation of oblique RSFs (“aorsf”) to the previous implementation, as well as other RSFs, GBMs, penalized CPH, and ANNs. They compare methods on 21 datasets, including low- and high-dimensional settings. Results are analysed using post-hoc Bayesian ROPE and evaluated using the Index of Prediction Accuracy (IPA) [54] based on the ISBS and time-dependent C-index [15]. They present results relative to the method performing best in their benchmark (“aorsf-fast” for both measures), with GBMs among the lowest-performing methods. Only minimal tuning was conducted.

A.1 Surveys of Survival Models

The final class of papers do not perform empirical benchmark experiments but instead survey/review available survival models. These are therefore only discussed very briefly. Ohno-Machado (2001) [77] provide an overview to models available for survival analysis from non-parametric estimators and classical models to neural networks. The review highlights useful applications of the models and their respective limitations. In particular their Table 1 clearly states advantages and disadvantages of Cox models versus ANNs. Patel *et al.* (2006) [78] compare proportional hazards and accelerated failure time models. This comparison is primarily theoretical and based on model properties, no analytical comparison with measures is provided though comparisons of predicted median survival times are compared to those from a Kaplan-Meier estimator. The authors conclude that AFT models should be considered more often due to simpler interpretation. Wang *et al.* (2019) [101] provide a review of survival analysis models and measures that is strongly recommended here as a precise and comprehensive introduction to the field of machine learning in survival analysis. The authors provide strong arguments for comparing classical models against one another, though this is not extended to the machine learning setting. More mathematical detail is provided to the classical setting however a clear and detailed overview is still provided for machine learning models. Some attention is also given to the more complex cases of competing risks and multiple events. Lee and Lim (2019) [67] provide a short but concise overview to survival analysis models with an emphasis on genetic data and implementation in R. Their review covers classical models, penalisation, and many machine learning models. They provide a clear, practical illustration (but no full benchmark experiment) comparing the models on a real dataset against Harrell’s C. No model outperforms CPH. Wiegrebbe *et al.* (2024) [102] provide a comprehensive overview of deep learning methods for survival analysis and compare methods based on their capabilities regarding various common challenges in survival analysis such as time-varying features, competing events, different censoring types, dimensionality,

modality, and interpretability. Their extensive comparison is also available as a web-based interactive table.

B Datasets

Table 2 lists datasets included in the benchmark along with their sources and common descriptive statistics. Section B lists the licenses declared by the source packages from which the dataset is taken.

Table 2: Datasets used in benchmark experiment.

Dataset ¹	Cens % ²	n_C^3	n_D^4	n^5	p^6	n_E^7	Package ⁸
aids.id [22]	60	1	4	467	5	188	<i>JM</i> [81]
aids2 [100]	38	1	3	2814	4	1733	<i>MASS</i> [100]
bladder0 [93]	48	0	3	397	3	206	<i>frailtyHL</i> [39]
CarpenterFdaData [23]	36	15	11	408	26	262	<i>simPH</i>
channing [58]	62	1	1	458	2	176	<i>KMsurv</i>
child [19]	79	1	3	26574	4	5616	<i>eha</i> [19]
colrec [72]	17	3	2	5578	13	4602	<i>relsurv</i> [72]
cost [50]	22	3	10	518	13	404	<i>pec</i> [73]
dataFTR [95]	86	0	2	2206	2	300	<i>RISCA</i> [32]
dataSTR [95]	82	0	4	546	4	101	<i>RISCA</i>
e1684 [57]	31	1	2	284	3	196	<i>smcure</i> [21]
flchain [29]	72	4	3	7871	7	1082	<i>survival</i>
gbsg [56]	43	3	4	2232	7	1267	<i>pycox</i> [61]
grace [45]	68	4	2	1000	6	324	<i>mlr3proba</i> [89]
hdfail [75]	94	1	4	52422	5	2885	<i>frailtySurv</i> [75]
kidtran [58]	84	1	3	863	4	140	<i>KMsurv</i>
liver [4]	40	1	1	488	2	292	<i>joineR</i> [103]
lung [70]	28	5	3	167	8	120	<i>survival</i>
metabric [56]	42	5	4	1903	9	1103	<i>pycox</i>
mgus [62]	6	6	1	176	7	165	<i>survival</i>
nafld1 [3]	92	4	1	12446	5	1018	<i>survival</i>
nwtco [18]	86	1	2	4028	3	571	<i>survival</i>
ova [98]	26	1	4	358	5	266	<i>dynpred</i>
patient [11]	79	2	5	1985	7	416	<i>pamntools</i> [10]
rdata [72]	47	1	3	1040	4	547	<i>relsurv</i> [72]
std [58]	60	3	18	877	21	347	<i>KMsurv</i>
support [56]	32	10	4	8873	14	2705	<i>pycox</i>
tumor [10]	52	1	6	776	7	375	<i>pamntools</i>
uis [44]	19	7	5	575	12	464	<i>quantreg</i> [59]
veteran [51]	7	3	3	137	6	128	<i>survival</i>
wbc1 [94]	43	2	0	190	4	109	<i>dynpred</i>
whas [45]	48	3	6	481	9	249	<i>mlr3proba</i>

1. Dataset ID and citation.

2. Proportion of censoring in the (modified) dataset, rounded to nearest percentage point.

3-4. Number of continuous and discrete features respectively before recoding.

5-6. Total number of observations and features respectively after alterations described above but before sub-sampling.

7. Number of observed events in dataset.

8. Package in which the dataset is included.

Table 3: Software licenses of the R and Python packages used as sources for datasets in this benchmark.

License	Packages
BSD-2	pycox
GPL	relnsurv
GPL (≥ 2)	jm, eha, pec, RISCA, dynpred, quantreg
GPL (≥ 3)	KMsurv
GPL-2	smcure
GPL-2 GPL-3	nnet
GPL-3	simPH, joineR
LGPL (≥ 2)	survival
LGPL-2	frailtySurv
LGPL-3	mlr3proba
MIT	pammtools
Unlimited	frailtyHL

C Models

Table 4 lists all the compared models, along with their software, version numbers, and prediction types. The horizontal lines separate the models into different groups: Baseline, Classical, Random Survival Forests (RSFs), Gradient Boosting Machines (GBMs) and Support Vector Machines (SVMs).

Table 4: Models used for benchmarking with associated packages and prediction types.

Model Information			Prediction Types		
Model Name ¹	Learner ²	Package ³	distr ⁴	crank ⁵	lp ⁶
Kaplan-Meier (KM) [52]	kaplan	survival v3.5.7	✓	ExpMort	×
Nelson-Aalen (NA) [1]	nelson	survival v3.5.7	✓	ExpMort	×
Akritas Estimator (AK) [2]	akritas	survivalmodels v0.1.18	✓	ExpMort	×
Cox PH (CPH) [26]	coxph	survival v4.1.8	✓ (Breslow)	lp	✓
CV Regularized CPH (GLMN) [84]	cv_glmnet	glmnet v4.1-8	✓ (Breslow)	lp	✓
Penalized (Pen) [36]	penalized	penalized v0.9.52	✓ (Breslow)	ExpMort	×
Parametric (AFT) [51]	parametric	survival v3.5.7	✓ (AFT)	lp	✓
Flexible Splines (Flex) [82]	flexible	flexsurv v2.2.2	✓	lp	✓
Random Survival Forest (RFSRC) [48]	rfsrc	randomForestSRC v3.2.2	✓	ExpMort	×
Random Survival Forest (RAN) [48, 105]	ranger	ranger v0.16.0	✓	ExpMort	×
Conditional Inference Forest (CIF) [46]	cforest	partykit v1.2.20	✓	ExpMort	×
Oblique Random Survival Forest (ORSF) [49]	orsf	aorsf v0.1.2	✓	ExpMort	×
Relative Risk Tree (RRT) [17]	rpart	rpart v4.1.23	(KM, PH)	✓	×
Model-Based Boosting (MBST) [20]	mboost	mboost v.2.9.9	✓ (Breslow)	lp	✓
CoxBoost (CoxB) [12]	cv_coxboost	CoxBoost v1.5	✓ (Breslow)	lp	✓
XGBoost (XGB Cox) [25]	xgboost	xgboost v1.7.6.1	✓ (Breslow)	lp	✓
XGBoost (XGB AFT) [7]	xgboost	xgboost v1.7.6.1	(KM, AFT)	lp	✓
SSVM-Hybrid (SSVM) [97]	svm	survivalsvm v0.0.5	(KM, PH)	✓	×

1. Identifier for the algorithm. Model abbreviations in parentheses are used in results.
2. Learner ID in *m3r*. Most learners are implemented in *m3r3extralearners* v0.7.1-9000 with KM, NA, CPH residing in *m3r3proba* v0.6.0
3. Package in which the learner is implemented with version used in experiment.
4. *distr* predict type in *m3r3proba* is the probabilistic prediction. A check (✓) represents the distribution being predicted directly by the package, in some cases using the Breslow estimator. The KM estimator may otherwise be used to estimate a baseline distribution from which the predicted distribution is composed with either PH or AFT forms, unless the Breslow estimator is used directly to obtain distribution predictions from the linear predictors. Notably for MBST, the availability of *distr* depends on the *family* parameter.
5. *crank* predict type in *m3r3proba* is the continuous ranking prediction. A check (✓) represents the ranking being predicted directly by the package. ‘ExpMort’ stands for expected mortality and is a risk score composed from the predicted survival distribution (*distr*). ‘lp’ represents the ranking being identical to the predicted linear predictor.
6. *lp* predict type in *m3r3proba* is the linear predictor prediction. A check (✓) represents the linear predictor being predicted directly by the package whereas a cross (×) means the prediction is not available (and cannot be composed).

D Model Configuration

Table 5 shows hyperparameter search spaces and non-default parameter values as well as common pre-processing requirements.

Table 5: Hyper-parameter search-spaces for tuning and non-default configurations for models.

Model	Hyper-parameters ¹	Values ²	Standardize ³	Encode ⁴
KM	-	-	×	×
NA	-	-	×	×
AK	lambda	[0, 1]	×	×
CPH	-	-	×	×
GLMN	alpha	[0, 1]	×	✓
Pen	lambda1	$2^{[-10,10]}$	×	×
	lambda2	$2^{[-10,10]}$		
AFT	dist	{weibull, lognormal, loglogistic}	×	×
Flex	k	{1, ..., 10}	×	×
	splitrule	{logrank, bs.gradient}		
RFSRC	ntrree	1000	×	×
	mtry	{1, ..., p}		
	nodesize	{1, ..., 50}		
	samptype	{swr, swor}		
	sampsiz	[0, 1]		
RAN	splitrule	{logrank,C,maxstat}	×	×
	num.trees	1000		
	mtry	{1, ..., p}		
	min.node.size	{1, ..., 50}		
	replace	{TRUE, FALSE}		
CIF	fraction	[0, 1]	×	×
	ntrree	1000		
	mtry	{1, ..., p}		
	minsplit	{1, ..., 50}		
	mincriterion	[0, 1]		
	replace	{TRUE, FALSE}		
fraction	[0, 1]			

Continued on next page...

Table 5: (continued)

Model	Hyper-parameters ¹	Values ²	Standardize ³	Encode ⁴
ORSF	control_type	fast		
	n_tree	1000		
	mtry	{1, ..., p}	×	×
	leaf_min_events	{5, ..., 50}		
	min_obs_to_split_node	min_events_to_split_node + 5		
RRT	alpha	(0, 1)		
	minbucket	{5, ..., 50}	×	×
MBST	family	{gehan, cindex, coxph, weibull}		
	mstop	{10, ..., 5000}	×	×
	nu	(0, 0.1]		
	baselearner	{bols, btree}		
CoxB	penalty	optimCoxBoostPenalty		
	maxstepno	5000	×	✓
	K	3		
XGBCox	objective	survival:cox		
	tree_method	hist		
	booster	gbtree		
	max_depth	{1, ..., 20}		
	subsample	[0, 1]	×	✓
	colsample_bytree	(0, 1]		
	nrounds	{10, ..., 5000}		
	eta	[10 ⁻⁵ , 10 ⁵]		
grow_policy	{depthwise, lossguide}			

Continued on next page...

Table 5: (continued)

Model	Hyper-parameters ¹	Values ²	Standardize ³	Encode ⁴
XGBAFT	objective	survival:aft		
	tree_method	hist		
	booster	gbtree		
	max_depth	{1,...,20}		
	subsample	[0, 1]		
	colsample_bytree	(0, 1]	×	✓
	nrounds	{10,...,5000}		
	eta	[10 ⁻⁵ , 10 ⁵]		
	grow_policy	{depthwise, lossguide}		
	aft_loss_distribution	{normal, logistic, extreme}		
aft_loss_distribution_scale	[0.5, 2]			
SSVM ⁵	type	hybrid		
	diff.meth	makediff3		
	gamma.mu	([2 ⁻¹⁰ , 2 ¹⁰], [2 ⁻¹⁰ , 2 ¹⁰])	✓	✓
	kernel	{lin_kernel, rbf_kernel, add_kernel}		
	kernel.pars	[2 ⁻⁵ , 2 ⁵]		

¹ Hyper-parameters for model tuning. The choice of hyper-parameters are largely informed by recommendations from the model author and subsequent papers exploring optimization. A '-' indicates no tuning is performed.

² Value ranges for the respective hyper-parameters to tune over. Omitted parameters use the package defaults.

³ Pre-processing of covariates by scaling to unit variance and centering to zero mean. A check (✓) indicates this step is performed before training the model, and a cross (×) if not.

⁴ Pre-processing of covariates by treatment encoding with `model.matrix`. A check (✓) indicates this step is performed before training the model, and a cross (×) if not.

E Implementation, Reproducibility, Accessibility

Platform All experiments were conducted on R 4.2.2 (2022-10-31) – “Innocent and Trusting” on the Beartooth Computing Environment [8].

Reproducibility and Accessibility Seeds were set with L’Ecuyer’s random number generator [63] to ensure reproducible results. All code required to run the experiments, as well as the results, are freely available in a public GitHub repository (https://github.com/slds-lmu/paper_2023_survival_benchmark). Software, packages, and version numbers that were utilized to conduct benchmark and analysis are listed below. We also employ the *renv* R package [96] to record and restore package dependencies to further aid reproducibility.

Packages All code is implemented using R v4.2.2 [79] and the experiment was run with *batch-tools* [64]. Models are implemented in *mlr3proba* v0.6.0 [89], *mlr3extralearners* v0.7.1-9000 [88] and *survivalmodels* v0.1.18 [85]. Tuning is implemented in *mlr3tuning* v0.19.2 [66]. Benchmarking functionality is implemented in *mlr3* v0.17.2 [65]. Benchmark analysis is implemented in *mlr3benchmark* v0.1.6 [87]. Compositions and pre-processing steps are implemented through *mlr3pipelines* v0.5.0-2 [13]. The packages and versions for the tested models are given in Table 4, all measures are implemented in *mlr3proba*.

F Results

The following plots show boxplots of the respective evaluation measure across the five outer re-sampling folds. Results are also available following links on the GitHub repository at https://github.com/slds-lmu/paper_2023_survival_benchmark.

As described in Section 4, we provide aggregated results based on average ranks across all datasets as boxplots with two versions for discrimination measures and three different scaling options for scoring rules:

1. Raw scores as produced by the evaluation measures, e.g. “ISBS”,
2. Explained Residual Variation (ERV) scores, e.g. “ISBS [ERV]”,
3. Scaled scores, e.g. “ISBS [Scaled]”, which scales raw scores such that 0 is the score achieved by KM and 1 is the score achieved by the best-performing model within the given combination of dataset and tuning- and evaluation measure.

Since the discrimination measures Harrell’s C and Uno’s C are already scaled from 0 (worst) to 1 (best) with KM achieving a score of 0.5 by design, the ERV option is omitted and only options 1) and 3) are presented.

F.1 Discrimination Measures

F.1.1 Raw Scores

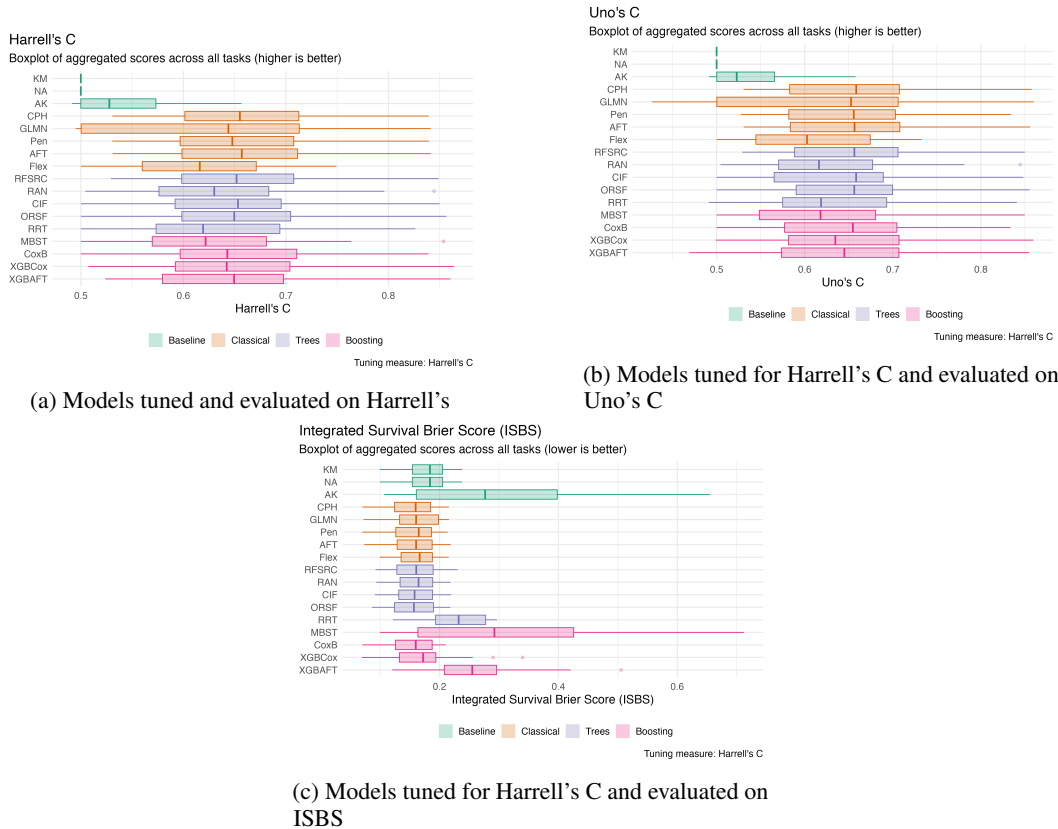
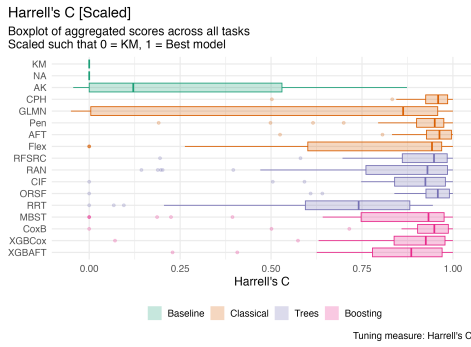
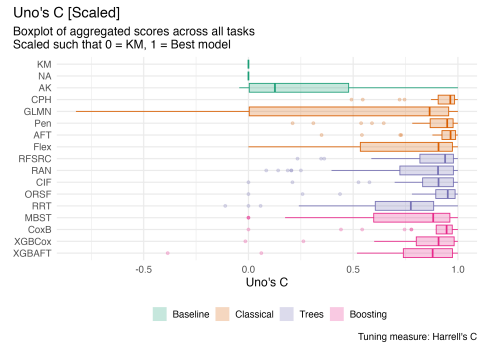


Figure 3: Boxplots of raw evaluation scores using discrimination measures for tuning (Harrell’s C) and using discrimination measures and ISBS for evaluation.

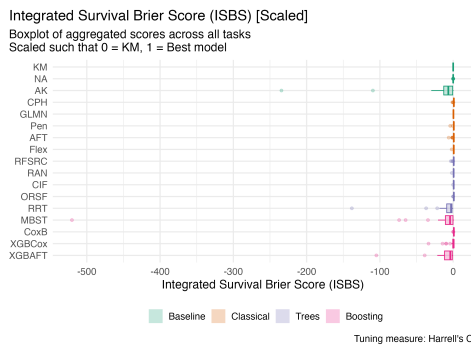
F.1.2 Scaled



(a) Models tuned and evaluated on Harrell's C



(b) Models tuned for Harrell's C and evaluated on Uno's C

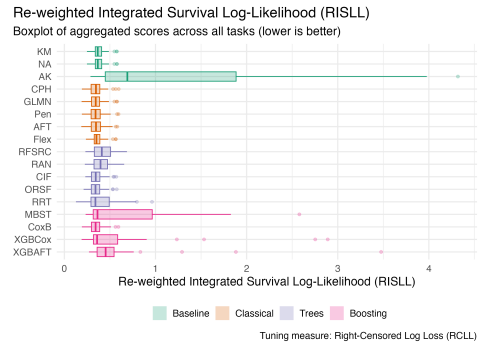
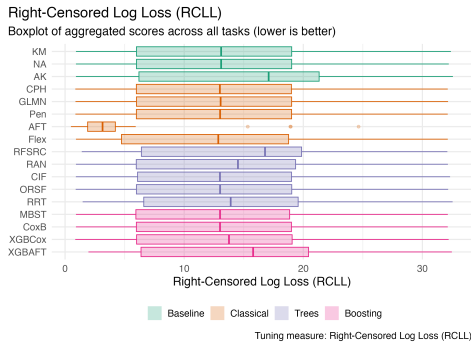


(c) Models tuned for Harrell's C and evaluated on ISBS

Figure 4: Boxplots of scaled evaluation scores using discrimination measures for tuning (Harrell's C) and using discrimination measures and ISBS for evaluation.

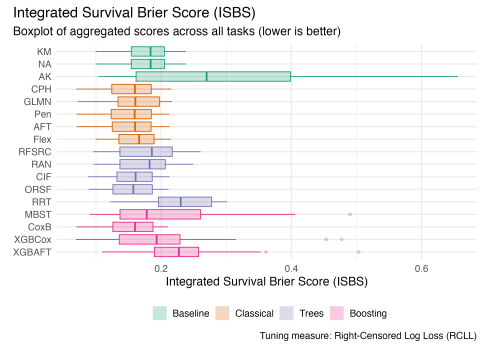
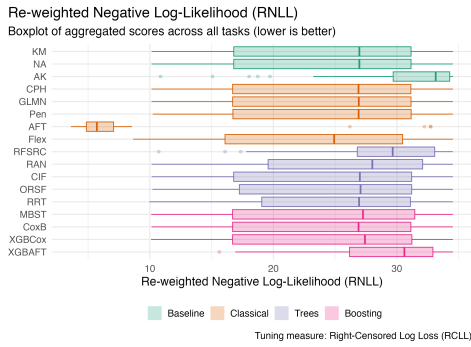
F.2 Scoring Rules

F.2.1 Raw Scores



(a) Scores for models tuned and evaluated on RCLL

(b) Scores for models tuned for RCLL and evaluated on RISLL

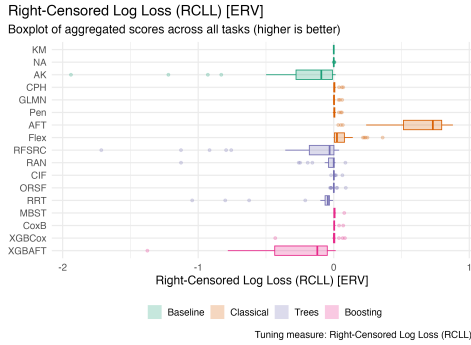


(c) Scores for models tuned for RCLL and evaluated on RNLL

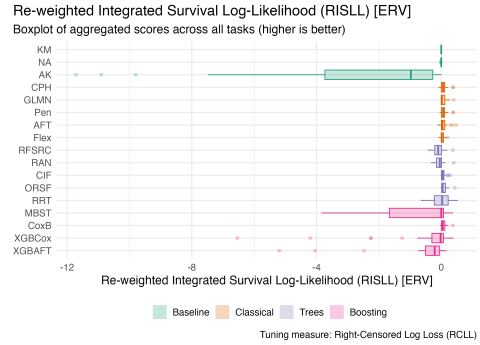
(d) Scores for models tuned for RCLL and evaluated on ISBS

Figure 5: Boxplots of raw evaluation scores using scoring rules for tuning (RCLL) and evaluation.

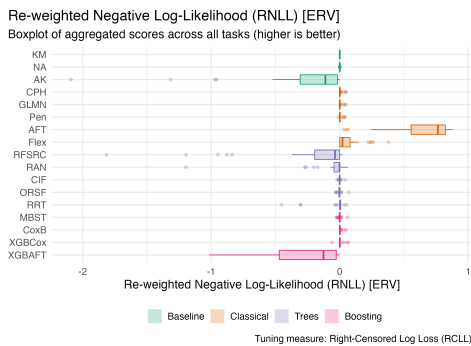
F.2.2 ERV



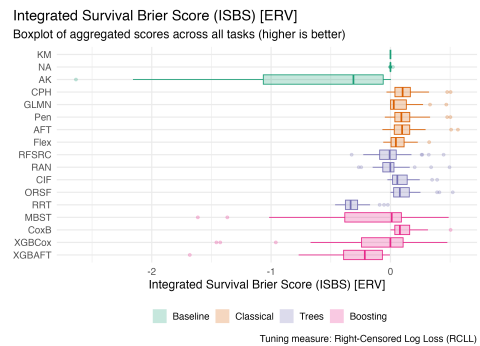
(a) ERV scores for models tuned and evaluated on RCLL



(b) ERV scores for models tuned for RCLL and evaluated on RISLL



(c) ERV scores for models tuned for RCLL and evaluated on RNLL



(d) ERV scores for models tuned for RCLL and evaluated on ISBS

Figure 6: Boxplots of ERV evaluation scores using scoring rules for tuning (RCLL) and evaluation.

F.3 Errors

During the many computational steps performed during this benchmark, software errors are inevitably bound to happen. As we noted in Section 3, we impute missing evaluation scores in resampling folds using the score of the Kaplan-Meier estimator, which affects the results presented in this paper.

Table 6 counts the number of errors encountered for each model, dataset, and tuning measures uses in the benchmark per outer resampling fold (up to five). We note that particularly the tasks hdfail and child have caused the majority of the runtime- or memory-related errors here due to their large sample sizes and number of unique time points.

Table 6: Number of errors per outer resampling iteration (up to five), separated by model, dataset, and tuning measure.

Model	Dataset	Harrell's C	RCLL	Total Errors
AK	child	4	4	8
AK	hdfail	5	5	10
GLMN	hdfail	0	1	1
Pen	child	0	2	2
Pen	hdfail	4	0	4
RFSRC	child	0	1	1
RFSRC	hdfail	1	1	2
RFSRC	nafld1	4	5	9
RAN	child	4	5	9
RAN	colrec	3	5	8
RAN	flchain	0	1	1
RAN	hdfail	4	5	9
RAN	metabric	2	2	4
RAN	nafld1	4	5	9
RAN	support	4	5	9
CIF	colrec	1	0	1
CIF	hdfail	2	2	4
CIF	nafld1	5	5	10
CIF	support	2	2	4
ORSF	nafld1	1	5	6
RRT	hdfail	0	2	2
MBST	aids2	0	3	3
MBST	child	5	5	10
MBST	dataFTR	1	1	2
MBST	flchain	3	2	5
MBST	gbsg	1	0	1
MBST	hdfail	5	5	10
MBST	metabric	1	0	1
MBST	nafld1	5	5	10
MBST	nwtco	4	5	9
MBST	support	4	1	5
CoxB	hdfail	5	5	10
XGBAFT	hdfail	0	1	1

G Results per Dataset

For completeness, we display boxplots of the evaluation scores across the outer resampling iterations per dataset, learner, evaluation measure and tuning measure.

G.1 Discrimination

Harrell's C

Scores per dataset across outer resampling folds (higher is better)



Figure 7: Raw scores for Harrell's C across outer resampling folds per dataset, task, and evaluation measure for models tuned on Harrell's C

Uno's C

Scores per dataset across outer resampling folds (higher is better)



Figure 8: Raw scores for Uno's C across outer resampling folds per dataset, task, and evaluation measure for models tuned on Harrell's C

Integrated Survival Brier Score (ISBS)

Scores per dataset across outer resampling folds (lower is better)

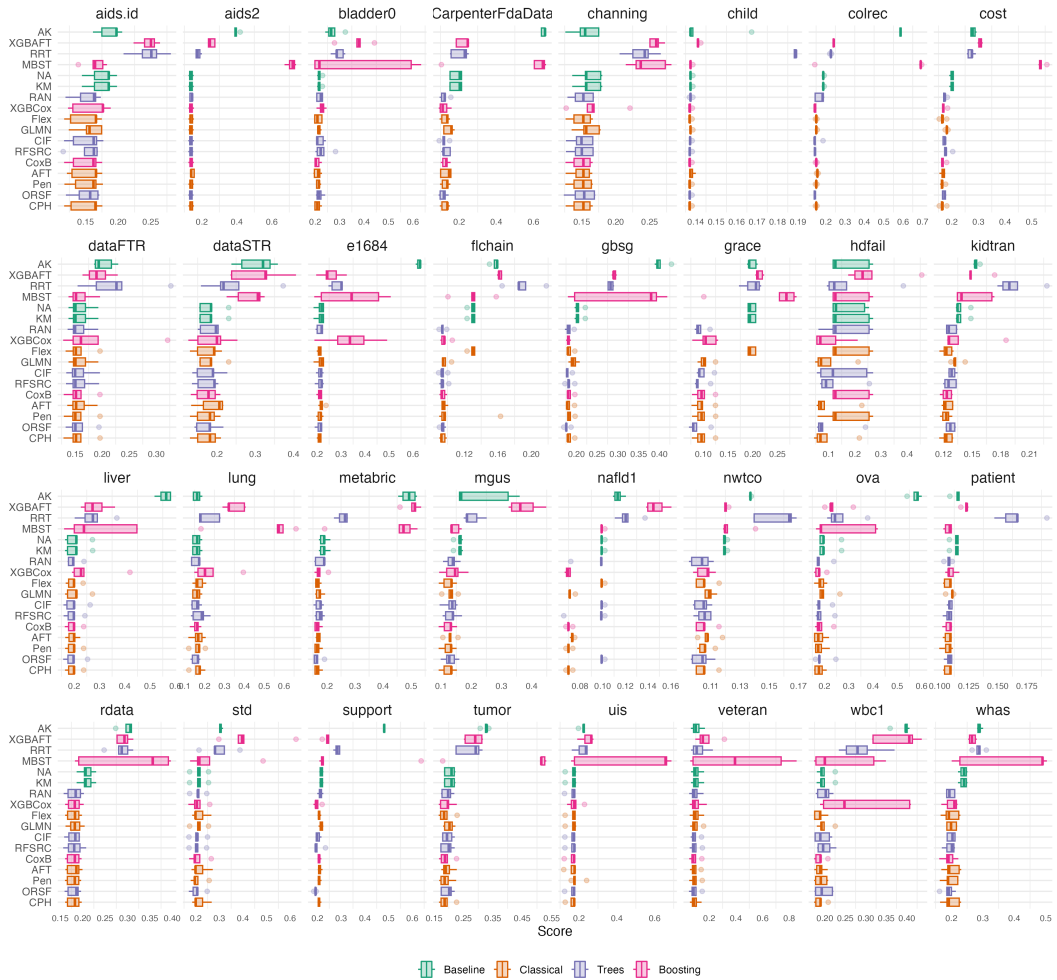


Figure 9: Raw scores for ISBS across outer resampling folds per dataset, task, and evaluation measure for models tuned on Harrell's C

G.2 Scoring Rules

Right-Censored Log Loss (RCLL)

Scores per dataset across outer resampling folds (lower is better)

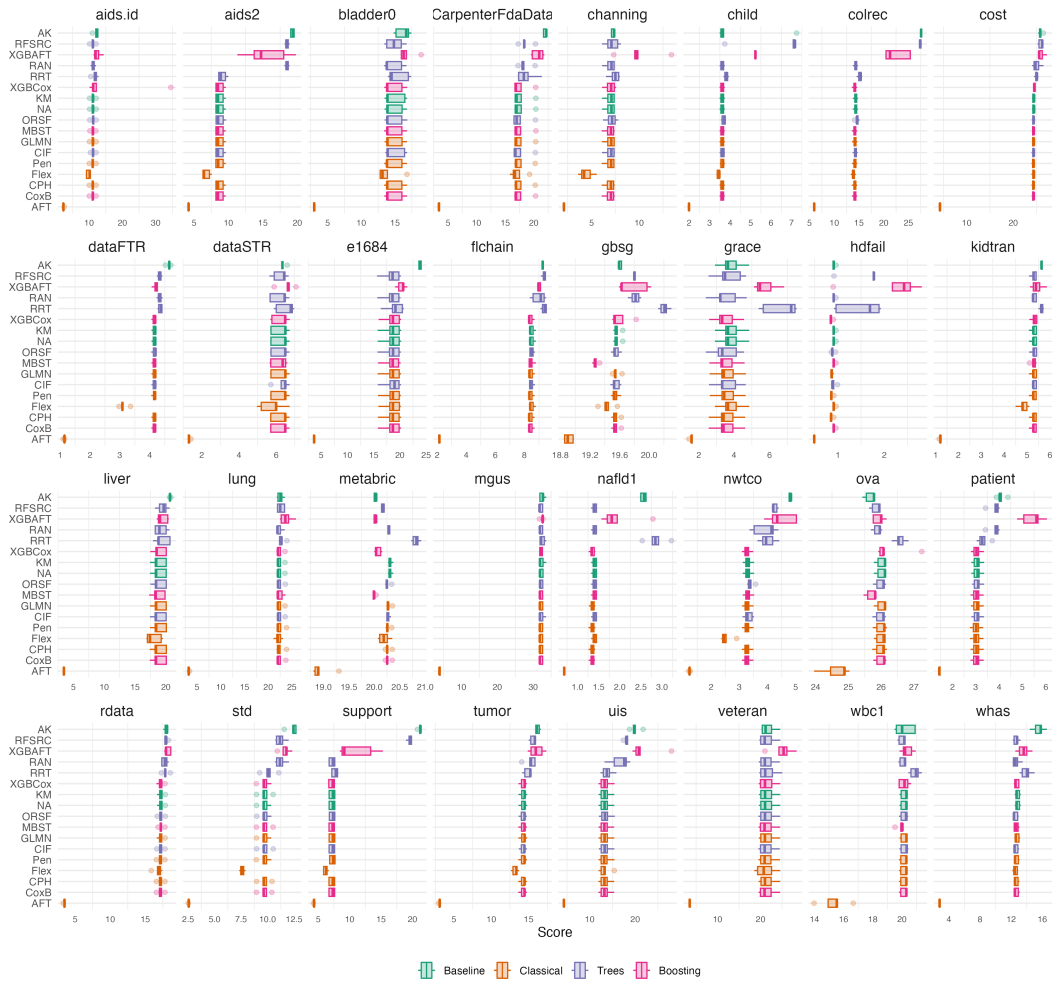
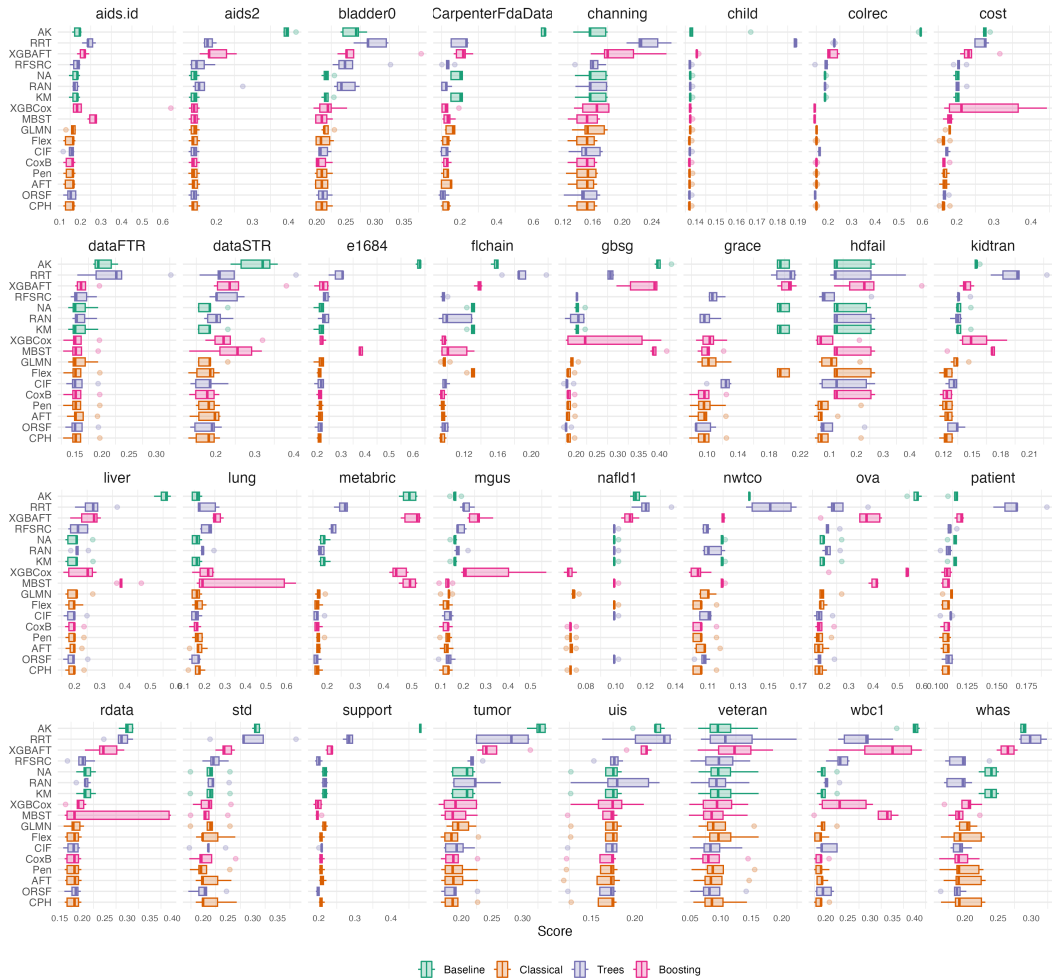


Figure 10: Raw scores for RCLL across outer resampling folds per dataset, task, and evaluation measure for models tuned on RCLL

Integrated Survival Brier Score (ISBS)

Scores per dataset across outer resampling folds (lower is better)



Tuning measure: Right-Censored Log Loss (RCLL)

Figure 11: Raw scores for ISBS across outer resampling folds per dataset, task, and evaluation measure for models tuned on RCLL

Re-weighted Negative Log-Likelihood (RNLL)

Scores per dataset across outer resampling folds (lower is better)

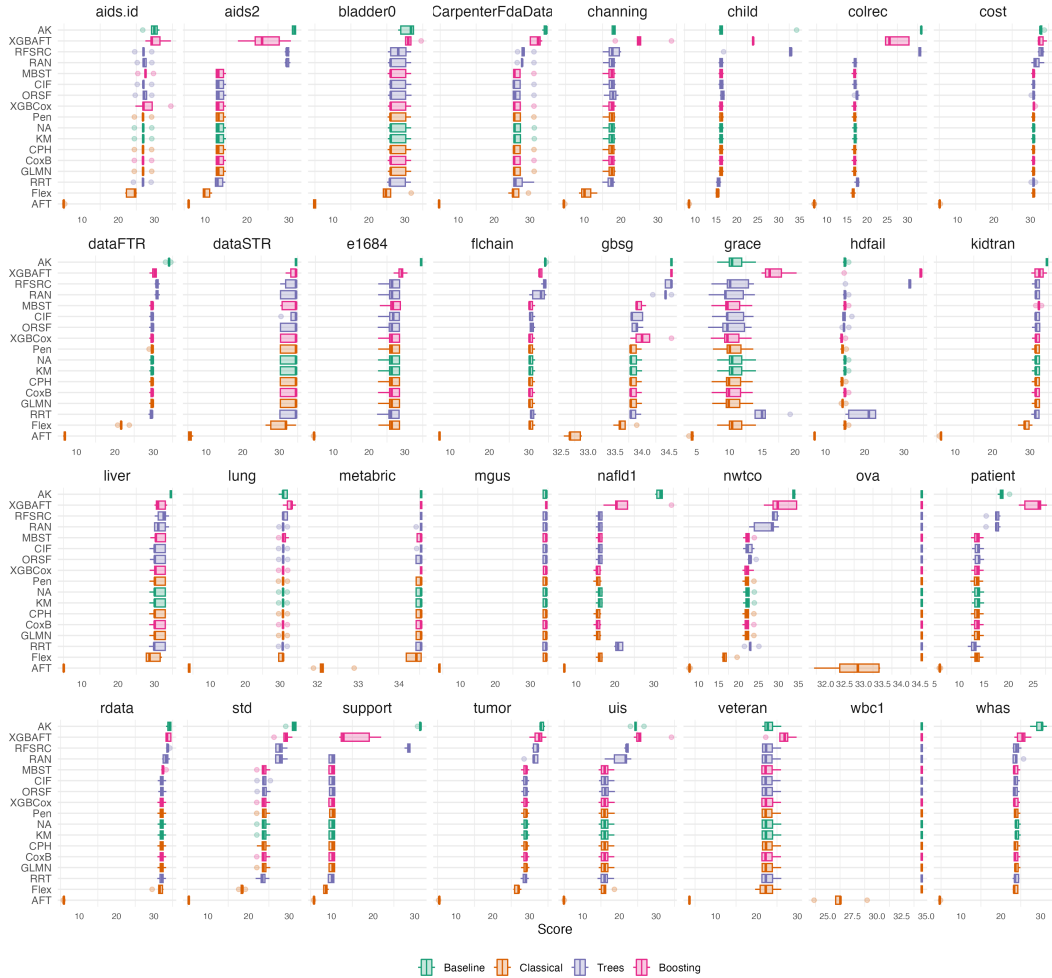


Figure 12: Raw scores for RNLL across outer resampling folds per dataset, task, and evaluation measure for models tuned on RCLL

Re-weighted Integrated Survival Log-Likelihood (RISLL)
 Scores per dataset across outer resampling folds (lower is better)

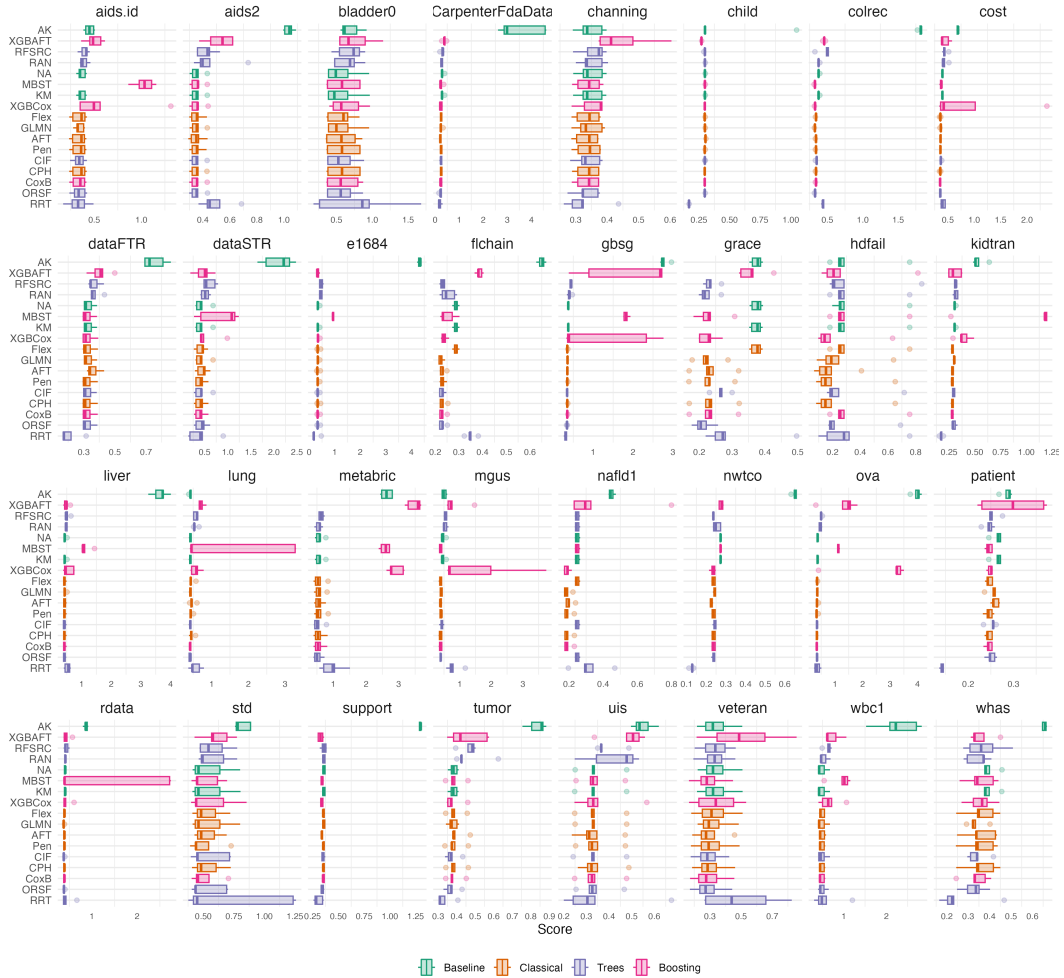


Figure 13: Raw scores for RISLL across outer resampling folds per dataset, task, and evaluation measure for models tuned on RCLL

G.3 Calibration

Both calibration measures here are presented only for models tuned on RCLL. We note that these measures may provide additional insights, yet consider both to be experimental in theory and implementation.

G.3.1 D-Calibration

D-Calibration considers a model well-calibrated if the p-value of the underlying test is greater than 0.05, which we display in form of a heatmap with X indicating a significant test result, indicating a model is not well-calibrated.

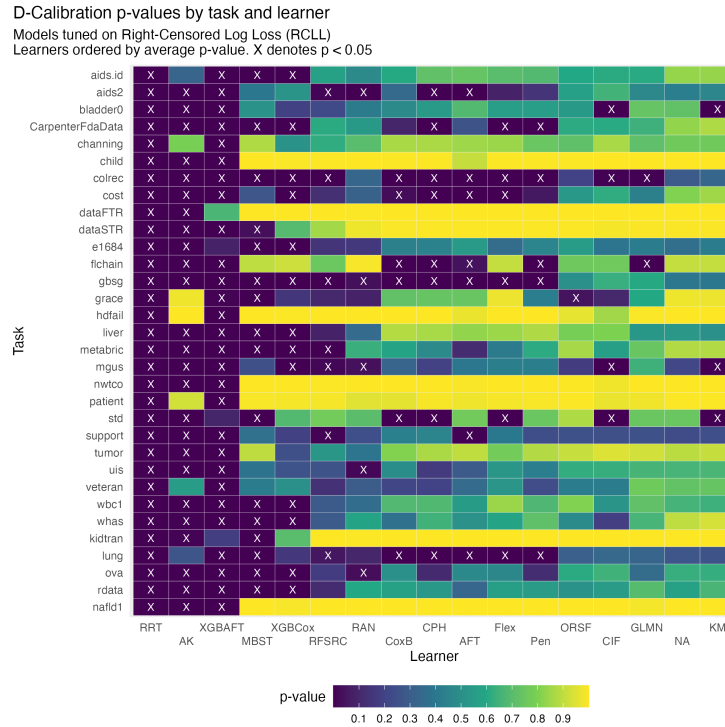


Figure 14: D-Calibration p-value heatmaps across all datasets and learners. ‘X’ indicates $p < 0.05$ while a non-significant result indicates good calibration.

G.3.2 Alpha-Calibration

Van Houwelingen's α relates predicted and observed hazards, with a value close to 1 implying a well-calibrated model.

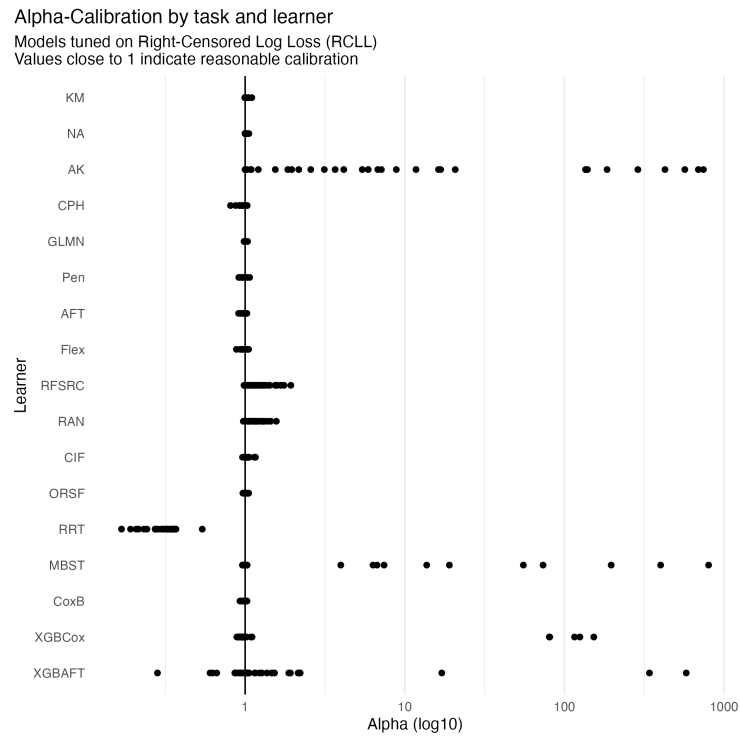


Figure 15: Calibration scores using van Houwelingen's α indicate good calibration when values are close to 1.