

# Development of an ensemble machine learning prognostic model to predict 60-day risk of major adverse cardiac events in adults with chest pain

Chris J. Kennedy PhD<sup>1, 2, \*</sup>, Dustin G. Mark MD<sup>1</sup>, Jie Huang PhD<sup>1</sup>, Mark J. van der Laan PhD<sup>3</sup>, Alan E. Hubbard PhD<sup>3</sup>, Mary E. Reed DrPH<sup>1</sup>

1. Kaiser Permanente Northern California, Oakland, CA
2. Department of Biomedical Informatics, Harvard Medical School, Boston, MA
3. Division of Biostatistics, UC Berkeley, Berkeley, CA

\* [chris\\_kennedy@hms.harvard.edu](mailto:chris_kennedy@hms.harvard.edu)

Version: March 12, 2021

## Abstract

**Background:** Chest pain is the second leading reason for emergency department (ED) visits and is commonly identified as a leading driver of low-value health care. Accurate identification of patients at low risk of major adverse cardiac events (MACE) is important to improve resource allocation and reduce over-treatment.

**Objectives:** We sought to assess machine learning (ML) methods and electronic health record (EHR) covariate collection for MACE prediction. We aimed to maximize the pool of low-risk patients that are accurately predicted to have less than 0.5% MACE risk and may be eligible for reduced testing.

**Population Studied:** 116,764 adult patients presenting with chest pain in the ED and evaluated for potential acute coronary syndrome (ACS). 60-day MACE rate was 1.9%.

**Methods:** We evaluated ML algorithms (lasso, splines, random forest, extreme gradient boosting, Bayesian additive regression trees) and SuperLearner stacked ensembling. We tuned ML hyperparameters through nested ensembling, and imputed missing values with generalized low-rank models (GLRM). We benchmarked performance to key biomarkers, validated clinical risk scores, decision trees, and logistic regression. We explained the models through variable importance ranking and accumulated local effect visualization.

**Results:** The best discrimination (area under the precision-recall [PR-AUC] and receiver operating characteristic [ROC-AUC] curves) was provided by SuperLearner ensembling (0.148, 0.867), followed by random forest (0.146, 0.862). Logistic regression (0.120, 0.842) and decision trees (0.094, 0.805) exhibited worse discrimination, as did risk scores [HEART (0.064, 0.765), EDACS (0.046, 0.733)] and biomarkers [serum troponin level (0.064, 0.708), electrocardiography (0.047, 0.686)]. The ensemble's risk estimates were miscalibrated by 0.2 percentage points. The ensemble accurately identified 50% of patients to be below a 0.5% 60-day MACE risk threshold. The most important predictors were age, peak troponin, HEART score, EDACS score, and electrocardiogram. GLRM imputation achieved 90% reduction in root mean-squared error compared to median-mode imputation.

**Conclusion:** Use of ML algorithms, combined with broad predictor sets, improved MACE risk prediction compared to simpler alternatives, while providing calibrated

predictions and interpretability. Standard risk scores may neglect important health information available in other characteristics and combined in nuanced ways via ML.

---

**Keywords:** chest pain, clinical predictive model, prognostic modeling, interpretable machine learning, ensemble learning, variable importance, accumulated local effects, generalized low-rank models

The omission of prediction from the major goals of basic medical science has impoverished the intellectual content of clinical work, since a modern clinician's main challenge in the care of patients is to make predictions.

Alvan Feinstein, 1983

## 1 Introduction

Chest pain is the second leading reason for emergency department visits (Rui et al. 2016) and is commonly identified as a leading driver of low-value health care. Workup protocols in patients with chest pain are designed to diagnose the potential for major adverse cardiac events (MACE). Missed diagnoses of MACE can be cause for medico-legal action, which may encourage conservative testing without health benefit. Accurate identification of patients at low risk of MACE is important to improve resource allocation and reduce overtreatment (Amsterdam et al. 2010). Risk scores aim to identify patients eligible for early discharge, avoiding additional stress testing and cardiac imaging that is unlikely to be of benefit (Greenslade et al. 2018). The primary biomarkers used for initial triage are elevated cardiac troponin, a sensitive marker of cardiac injury measured serially, and repeated electrocardiograms.

Previous work has focused on the development and validation of additive risk scores as decision aids for risk stratification. Such risk scores examine a small number of biomarkers and demographics, summarize those predictors into qualitative levels, and use a weighted sum to allocate patients into risk categories. Standard risk scores are HEART (History, ECG, Age, Risk factors and Troponin) and EDACS (Emergency Department Assessment of Chest Pain Score - Than, Flaws, et al. 2014). HEART is most commonly used in North America, although EDACS has similar performance characteristics (Mark et al. 2018). Effective risk scores will stratify patients across risk levels such that the qualitative "low risk" group will have sufficiently low risk of short-term MACE that those patients can be discharged without additional workup. An ineffective or ill-calibrated risk score would underestimate the risk in the "low risk" group and lead to an overly optimistic early discharge policy that results in increased future MACE. But given multiple risk scores that are well-calibrated, scores with improved discrimination could theoretically result in a larger percentage of low-risk patients.

### 1.1 Background and Objectives

It remains debated whether machine learning methods can exhibit statistically and substantively significant benefits for risk prediction compared to logistic regression, decision trees, or additive risk scores (Goldstein, Navar, and Carter 2016; Goldstein, Navar, Pencina, et al. 2016). A recent meta-analysis, for example, did not find systematic benefit from machine learning in comparison to logistic regression (Christodoulou et al. 2019). Yet there is also optimism about the potential for artificial intelligence methods in medicine (He et al. 2019) in general, as well as cardiology specifically (Johnson et al. 2018).

Building on Mark et al. 2018, we sought to assess the performance of machine learning (ML) methods at predicting MACE among emergency department patients with chest pain. Could ML improve upon existing validated risk scores through a more complex integration of predictors that can better estimate MACE risk? To what extent is hyperparameter optimization necessary to achieve strong ML performance?

Our clinical objective was to maximize the pool of low-risk patients that are accurately predicted to have less than 0.5% MACE risk and may be eligible for reduced testing. The primary threshold of 0.5% risk has previously been identified as an acceptable risk by a majority of emergency physicians for early discharge (Than, Herbert, et al. 2013). Using a risk of 0.5% as the test threshold will inherently lead to a negative predictive value of greater than 99.5%, provided that the risk prediction is well-calibrated in the target population. We also examined secondary thresholds of 1.0% and 2.0%.

A reasonable assessment of ML performance could only be made in comparison to realistic alternative options. We compared ML performance to simpler indicators of risk: key biomarkers (troponin, electrocardiogram), validated clinical risk scores (History, ECG, Age, Risk factors and Troponin [HEART] and Emergency Department Assessment of Chest pain Score [EDACS]), decision trees, and logistic regression.

If machine learning can demonstrate improved discriminative performance compared to logistic regression and related methods, along with appropriate calibration, its next hurdle for adoption is to provide interpretability. Clinicians may be willing to forgo maximum predictive accuracy for the sake of understanding how individual predictors influence the output of the algorithm. With analytical effort it may be possible to provide sufficient interpretability for clinicians to accept the complication of machine learning and the benefit of the (potentially) improved predictive accuracy. To facilitate interpretation, we explained the models through prediction-based variable importance ranking and accumulated local effect visualization. If simpler algorithms remain preferred, the ML results can at least approximately the best achievable performance, and so serve as benchmark standards when considering more restrictive algorithms.

Certain analytical characteristics would be important to arrange in order for ML to potentially improve upon simpler options. First, it was important to extract a broader set of granular predictor variables than were used by existing scores. Extensive predictor sets give ML the potential to capture interactions and nonlinear relationships that are missed by linear or additive approaches, perhaps relevant only to certain subgroups of patients. Further, ML may statistically identify novel predictors that have been missed by existing scores or the broader literature, or whose predictive impact was too small, in too complex a form, or underrepresented in terms of sample size to be detected by non-ML methods. The expansion of electronic health records (EHRs) also makes broader covariate collection more feasible and relevant than was possible prior to EHRs, while also facilitating more granular measurement of variables (E. H. Kennedy et al. 2013).

It is also important for variables be measured on a fine-grained scale, which gives ML the opportunity to detect novel cut-points or thresholds that improve performance. Variables should be kept as their original continuous measurements rather than dichotomized or discretized into qualitative levels (Senn 2005). For example, a predictor such as body mass index (BMI) loses substantial information when it is dichotomized into an indicator of high-BMI or the absence of high-BMI. A single threshold chosen for that dichotomization may not be optimal for certain subgroups or regions of risk. One of the benefits of ML is that it can identify thresholds in a data-adaptive way, allowing it to better approximate unknown or ill-understood physiological processes. That said, very high cardinality variables can result in overfitting and slow down the training of certain machine learning algorithms, such as decision trees, that test each unique value as a possible subgroup splitting threshold. It may be beneficial to reduce the cardinality of granular continuous variables through histogram binning that scales with the dataset size, e.g.  $\text{sample size} / 1000$  or  $\log(\text{sample size}) * 10$ .

## 2 Data and Methods

### 2.1 Source of data

Our study was sourced from the electronic health record (EHR) of 21 emergency departments (EDs) within Kaiser Permanente Northern California, an integrated health care delivery system with over 1 million annual ED visits. patient visits to

### 2.2 Participants

All adult patients were retrospectively included if they had received cardiac troponin testing in the emergency department and either presented with a chief complaint of chest pain or chest discomfort, or whose ED physician had assigned them a primary or secondary ICD-coded diagnosis of chest pain. The later inclusion criterion is important because patients may complain of “anginal equivalents” (such as shortness of breath) in lieu of overt chest pain (Amsterdam et al. 2010). The initial inclusion pool had a 60-day MACE rate of 8.0%. Patients were excluded if they had a MACE diagnosis in the ED or within 30 days prior to ED visit, alternative non-ACS diagnoses at index visit (e.g. pneumonia, pneumothorax, or traumatic injury), could not be tracked due to lack of active health plan membership during the study (except in cases of death), or had a troponin I > 99th percentile upper limit of normal given the dominant predictive value of elevated troponin values for adverse outcomes in both patients with acute coronary syndromes and in the general population (Bonaca et al. 2010; De Lemos et al. 2010). Patients were excluded if their smoking status was unknown, which was viewed as a key marker of low-quality data. The final study cohort consisted of 116,764 patients with a 60-day MACE incidence of 1.88%. A fourth-generation troponin assay was used during the study period (AccuTnI+3, Beckman-Coulter, Brea, CA, USA).

### 2.3 Outcome

Our primary outcome was cumulative MACE incidence within 60 days of the index visit. We defined MACE as myocardial infarction, cardiogenic shock, cardiac arrest, or death.

### 2.4 Predictors

We used a total of 74 predictors sourced from the electronic health record, including vitals, labs, history, qualitative interpretation of ECG imaging, regular expression-based extraction of features from clinical notes, demographics, and missingness indicators (20). These predictors are detailed in Table ??.

### 2.5 Missing data

Missingness rates for each predictor are listed in Table ?. We created missingness indicators for each predictor, which marked the observations that were missing a value. Inclusion of missingness indicators often improves predictive performance (Agor et al. 2019). That matrix of missingness indicators was analyzed for perfect collinearity, and duplicate indicators were dropped.

Missing predictor values were imputed by factorizing the raw data matrix with generalized low-rank models (GLRM) (Schuler et al. 2016; Udell et al. 2016). GLRM is a generalization of principal component analysis and matrix completion methods and is designed for mixed type data frames that include continuous, categorical, ordinal, and binary variables. GLRM decomposes (factorizes) the original data frame into an  $X$  matrix of reduced components and  $Y$  matrix of archetypes, including possible penalty

terms that can induce sparsity (L1) or simply denoise (L2 or quadratic). Multiplying these two factor matrices reconstructs the original data frame, imputing any data entries with missing values. The method used for missingness provides few constraints on the resulting fit and also permits prediction from future data with missing values.

The GRLM hyperparameter settings were chosen through a grid search in which each model was trained on 75% of the data and evaluated on the remaining 25% for accuracy at reconstructing the original observed data matrix. Missingness indicators were not included in the GRLM imputation analysis. Our final GRLM settings were: 50 components, quadratic regularization on X with weight 4, and L1 regularization on Y with weight 24. Cells with missing data were then replaced with the reconstructed data matrix from GRLM using the optimal settings.<sup>1</sup>

GRLM imputation greatly increased the number of unique values (cardinality) for continuous variables, which would have a negative performance impact on tree-based algorithms that test every unique value for a potential split. To avoid that performance drop, we used penalized histogram binning to bin imputed predictors with high cardinality into up to 200 unique values (Rozenholc et al. 2010).

Multiple imputation was not necessary because our scientific goal was to characterize predictive performance for the unimputed outcome variable, rather than to estimate statistical parameters for covariates that were imputed, such as linear regression coefficients (Steyerberg 2009; Wang et al. 1992).

## 2.6 Prediction algorithms

Dozens if not hundreds of other prediction algorithms would be possible to evaluate, but computational time limitations forced us to choose a finite set with reasonable performance expectations. We chose well-known prediction algorithms that have shown strong performance in prior research, including both linear and decision tree-based estimation. The tree-based prediction algorithms were random forest (Breiman 2001), extreme gradient boosting (XGBoost) (Chen et al. 2016), and Bayesian additive regression trees (Chipman et al. 2010). The linear prediction algorithms were generalized additive models (T. J. Hastie et al. 1990) using thin plate splines (Wood 2003), and lasso (Tibshirani 1996).

Splines have shown competitive performance with tree-based algorithms in prior clinical prediction work due to their ability to identify non-linear, but smooth patterns (Austin 2007). The lasso algorithm (or its generalization the elastic net) is a helpful test of sparsity in the covariates, and a faster & more nuanced variable selection method than best subset or stepwise selection (T. Hastie et al. 2017). Better performance for lasso compared to logistic regression would indicate that feature selection could be helpful for other algorithms, while equal performance could indicate that the extraction of predictors from the EHR was overly restrictive and should be broadened.

## 2.7 Benchmarks

When evaluating complex algorithms it is important to contextualize their performance by comparing to simpler alternative approaches or benchmarks. If the benchmark algorithms can achieve similar performance then the extra complexity of the statistical machine learning algorithms may not be worthwhile. The improvement of a novel prediction method over standard benchmarks is known as the *skill* of the prediction method (Brier 1950; Murphy et al. 1977; F. Sanders 1963). In clinical prediction the primary alternatives to statistical machine learning are relatively inflexible fits, which

<sup>1</sup>Here X refers to the reduced components after GRLM transformation, and Y refers to the complementary matrix that transforms those components back to the original covariate space. It does not refer to the outcome variable.



include logistic regression, ordinary least squares, individual decision trees, and stratification on key clinical covariates. We tested each of these options, where key covariates were defined as peak troponin, qualitative ECG reading, EDACS score, and HEART score. As a complement to stratification on different subsets of key covariates, we also evaluated logistic regression and decision trees when restricted to these key covariates.

## 2.8 Stacked ensembling: SuperLearning

When comparing a variety of algorithms an initial choice is to use cross-validation to select the algorithm with the best out-of-sample performance. A more nuanced decision would be to consider a weighted average of multiple algorithms - creating a team of algorithms whose contribution to the prediction is based on optimizing out-of-sample performance on a certain statistic. That is the nature of stacked ensembles (Breiman 1996; Wolpert 1992), sometimes referred to as the Super Learner algorithm (van der Laan et al. 2007). Rather than restrict our prediction machine to a single algorithm, we create a weighted average across all tested algorithms, and estimated weights based on an optimization goal so that they minimize a chosen performance statistic on test data. We chose to optimize the Brier score (i.e. mean-squared error) in our ensemble, using convex weights based on a non-negative least squares meta-learner. Optimizing on Brier score includes a focus on both discrimination and calibration for the ensemble (Murphy et al. 1977). A convex combination of algorithm weights ensures that predictions fall within the convex hull of the constituent learners, while also inducing sparsity - i.e. algorithms can have zero weight.

## 2.9 Hyperparameter tuning

Prediction algorithms often have multiple hyperparameter settings that adjust the estimation procedure in different ways. Those hyperparameters are not estimated from the data, but rather must be specified a priori by the analyst. While software implementations will typically provide a default value for each hyperparameter, there is no reason to believe that the default values are effective for the current dataset. Customizing the hyperparameter configuration to the current dataset can allow the algorithms to adapt to the available sample size, number of predictor variables, measurement error in the predictors, sparsity in predictor relevance, and correlation structure of the predictors. Hyperparameters are often chosen by fitting the algorithm with different configurations and selecting the configuration that maximizes accuracy on held-out data, such as through cross-validation. The benefit of hyperparameter tuning is believed to vary by algorithm, which is referred to as the *tunability* of the algorithm (Probst, Boulesteix, et al. 2019). Random forest, for example, is believed to work well with default hyperparameters but also can benefit from hyperparameter tuning, particularly to reduce overfitting (Probst, Wright, et al. 2019; Segal et al. 2011).

Hyperparameter tuning is inherently a computationally intensive process, as it involves fitting the algorithms many different times, and varies based on the number of hyperparameters (dimensionality) as well as number of the unique values tested for each hyperparameter (resolution). Further complexity is involved if one considers that some hyperparameters may be more important than others for a given algorithm. Given the role of hyperparameters in modifying the performance of prediction algorithms, caution is warranted when generalizing algorithm performance characteristics from individual studies (e.g. algorithm X outperforms algorithm Y), particularly when hyperparameters are left at their default values and therefore are not customized to the given dataset.

For this work we adopted a hyperparameter tuning approach using *nested ensembling*. Much as using a weighted ensemble of different algorithms may be

preferable to selecting the single best-performing algorithm, using a weighted ensemble of hyperparameter settings for a given algorithm may yield improved performance compared to selecting a single set of hyperparameters. With that concept in mind we created small grids of hyperparameter configurations and estimated a SuperLearner ensemble for a given algorithm in which the ensemble weights selected the hyperparameter settings that maximized out-of-sample performance. This ensemble of hyperparameter settings could potentially rely on a single configuration due to the sparsity induced by the convex combination, or the optimization could distribute the weighting across multiple configurations if such a weighting improved performance over a single selected configuration. Another benefit of the nested ensembling is that it limits the number of learners that are analyzed in the outer SuperLearner ensemble, which can conserve power and mitigate overfitting in the meta-learning process (i.e. allocation of weights in the convex combination).

We used the ensemble hyperparameter tuning approach for random forests, xgboost, and individual decision trees. The random forest grid consisted of 9 configurations: minimum node size  $\in \{5, 20, 60\} \times$  covariates sampled  $\in \{4, 8, 16\}$ <sup>2</sup>. The xgboost grid consisted of 8 configurations: number of trees  $\in \{250, 1000\} \times$  maximum tree depth  $\in \{2, 4\} \times$  shrinkage  $\in \{0.05, 0.2\}$ . The decision tree grid consisted of 12 configurations: complexity parameter  $\in \{0, 0.01\} \times$  minimum split  $\in \{10, 20, 80\} \times$  maximum tree depth  $\in \{10, 30\}$ .

## 2.10 Evaluation

We evaluated alternative options for risk prediction based on their discrimination, calibration, and clinical utility. Nested cross-validation with 5 folds was used to conduct the discrimination and calibration analyses. While bootstrap estimation has been promoted for evaluation of clinical prediction models (Austin and Tu 2004; Steyerberg, Harrell Jr, et al. 2001), recent work has shown that the bootstrap can be biased for evaluating the performance of highly adaptive ML algorithms estimators such as random forests (Benkeser et al. 2019).

### 2.10.1 Discrimination

We chose area under the precision-recall curve (PR-AUC, also known as average precision) as our primary performance metric for evaluating discrimination, because it highlights performance differences that may be missed by ROC-AUC with imbalanced data (Cook 2007; Saito et al. 2015). We included area under the receiver operating characteristic curve (ROC-AUC or the concordance statistic) as our secondary performance metric, which remains highly popular and interpretable (Janssens et al. 2020). As an exploratory metric we also estimated the adjusted Brier score (index of prediction accuracy) which integrates discrimination and calibration into a single metric (Kattan et al. 2018). We visualized improvements in discriminative performance using density plots of the calibration slope (Steyerberg, Vickers, et al. 2010). We did not conduct a reclassification analysis due to recognized limitations (Hilden et al. 2014; Kerr et al. 2014; Leening et al. 2014; Pepe et al. 2015).

### 2.10.2 Calibration

Our clinical use case was centered on a risk threshold of 0.5% to classify patients as “low risk” in order to qualify for early discharge. Because of that scientific goal, it was especially important to compare the model’s predicted risks to the observed risks, i.e.

<sup>2</sup>The number of covariates sampled (i.e.  $m_{try}$ ) was based on the formula:  $\text{floor}(\{0.5, 1, 2\} \cdot \sqrt{p})$  where  $p$  is the total number of covariates.



its *calibration* (Lichtenstein et al. 1981) - also known as reliability (Brier 1950; Murphy et al. 1977) or external correspondence (Yates 1982). We assessed the calibration of predicted probabilities in two ways: 1) calibration curve visualization, 2) calculation of the index of prediction accuracy (IPA), a transformation of the Brier score (Kattan et al. 2018). We did not conduct a Hosmer-Lemeshow group-based calibration test due to its recognized limitations and recommendations against its use (Kramer et al. 2007; Van Calster et al. 2019).

### 2.10.3 Clinical utility

The planned clinical use of the prediction model was first to assess eligibility for early discharge among low-risk patients. Accurately estimating the risk of MACE for patients would allow those low-risk patients to be discharged and avoid additional unnecessary workup, freeing up resources (clinical attention, testing capacity, etc.) for higher risk patients. Low risk was generally defined as being below a 0.5% well-calibrated probability of MACE within 60 days, with less conservative thresholds of 1% and 2% as additional options.

Our model needed to balance two trade-offs: 1) false “negatives” in which a patient was identified as low-risk but whose true risk of MACE within 60 days was above the threshold, and 2) false “positives” in which patients were believed to be above the given threshold but whose true risk was less than the threshold. Errors in the first category have a greater cost than those in the second category, because there is a greater potential detriment to those patients who were discharged early but whose true risk exceeded the threshold. Patients incorrectly estimated to be above the risk threshold, but who are truly low risk, have comparatively minor costs of additional workup, use of clinical resources, and potential to be overtreated. Yet these possible errors are not quite the same as false negative or false positives typically used to assess predictive models: we care about the true, but unknown, risk rather than the observed outcome. Under this decision-making calculus a patient whose true risk is correctly predicted to be below the clinical threshold, and is therefore discharged without additional workup, but who ends up having a MACE would still have been managed appropriately.

This suggests that the absolute or squared error of the patient’s predicted risk versus true risk, particularly near the clinical threshold, would be reasonable loss functions to translate into clinical utility. Miscalibration near the clinical threshold needs to be avoided, whereas miscalibration away from the threshold does not affect the decision. As the expected value of that loss approaches zero we would see that the number of false negatives and false positives (in terms of risk above or below a threshold rather than the observed outcome) also approaches zero. We could target a specific threshold by focusing on patients on the incorrect side of the threshold and averaging the error in their risk prediction, possibly including differential weights for each side of the threshold to account for different costs to the patient. Such a “miscalibration-around-a-threshold” loss function might look as follows:

$$\text{loss}(Y_i, \hat{Y}_i | X_i) = \omega_1 \mathbf{1}(P_0(Y_i | X_i) < \tau) g(\hat{f}(X_i) - \tau) + \omega_2 \mathbf{1}(P_0(Y_i | X_i) > \tau) g(\tau - \hat{f}(X_i)) \quad (1)$$

where:

- $Y$  is the observed outcome and  $X$  is the set of predictors,
- $i$  indexes each patient in the sample,
- $P_0(Y_i | X_i)$  is the true risk of patient  $i$ ,

- $\hat{f}(X_i)$  is the predicted risk of patient  $i$  from a given estimator  $\hat{f}$ , 319
- $\tau$  is the clinical threshold (e.g. 0.5%), 320
- $g$  is a function such as the identity, squared value, or absolute value function, 321
- $\omega_1$  is the differential cost for low-risk patients who are kept for further workup, 322
- $\omega_2$  is the differential cost for high-risk patients who are incorrectly discharged early, 323  
324

We do not know the true risk for any patients, but we can estimate it within our sample by fitting a semi-parametric smooth function (e.g. lowess) to estimate the true probability of the outcome given the estimated predicted probability, equivalent to what is done during calibration analysis. 325  
326  
327  
328

If multiple decisions were to be made based on the estimated risk, we might sum this loss over each decision. Alternatively we might use a threshold-free loss function, such as: 329  
330  
331

$$\text{loss}(Y_i, \hat{Y}_i | X_i) = g(\hat{f}(X_i) - \tau) \quad (2)$$

In this work we focus on the threshold-free loss with absolute value as the transformation function  $g$ . 332  
333  
334

## 2.11 Interpretability 335

Beyond the statistical performance of a clinical prediction, it can be important to provide an explanation or overview of how a model generates its predictions. Interpretation is desirable first because it can provide evidence that the model is working as expected, which can improve the trustworthiness of its predictions for clinicians, patients, or collaborators. Interpretation may also lead to scientific insights about how predictors are related to the outcome, which could be conceptualized as causal pathways, data generating processes, or biological mechanisms. Interpretation can further inform the data export and cleaning processes, such as identifying extreme values, data entry errors, or outliers, or suggesting additional predictor variables to incorporate into the model. 336  
337  
338  
339  
340  
341  
342  
343  
344  
345

Methods of interpretation can be model-specific or model-agnostic. For models within the family of linear regression, one might provide the estimated beta coefficients for each predictor, along with their associated confidence intervals and p-values. Interpretation becomes less straightforward as models become more complex, such as with interaction or polynomial terms in a regression, random forest or boosted tree models with hundreds or thousands of non-linear decision trees, or splines in which ranges of a given predictor might have different coefficients. 346  
347  
348  
349  
350  
351  
352

In this work we focus on two complementary forms of model interpretability: *variable importance ranking* and *accumulated local effect plots*, as described below. 353  
354

### 2.11.1 Variable importance ranking 355

Prediction-oriented variable importance rankings order the predictor variables by their contribution to a model's prediction, providing evidence as to which predictors were relied upon the most by the algorithm. Such rankings could be used as a form of confirmatory analysis if a hypothesized ranking were created prior to data analysis, which could formally identify predictors that differed from their expected importance. 356  
357  
358  
359  
360

### 2.11.2 Accumulated local effects 361

It may also be helpful to understand how a model's prediction varies over the values of individual predictors, particularly continuous predictors with a wide range or large number of unique values. Partial dependence plots (PDPs) as proposed by Friedman (2001) are commonly used to provide this type of interpretability, but they can yield flawed results because they make a key unrealistic assumption that features are statistically independent of each other (Molnar 2020, p. 5.1.3). Accumulated local effect (ALE) plots are a recently developed method that avoids that limitation of PDPs, by counterfactually modifying observations that lie within a nearby kernel neighborhood of the current predictor's value of interest (Apley et al. 2019). Following the variable importance ranking, we visualize the contribution of high-importance continuous variables using accumulated local effect plots. 362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372

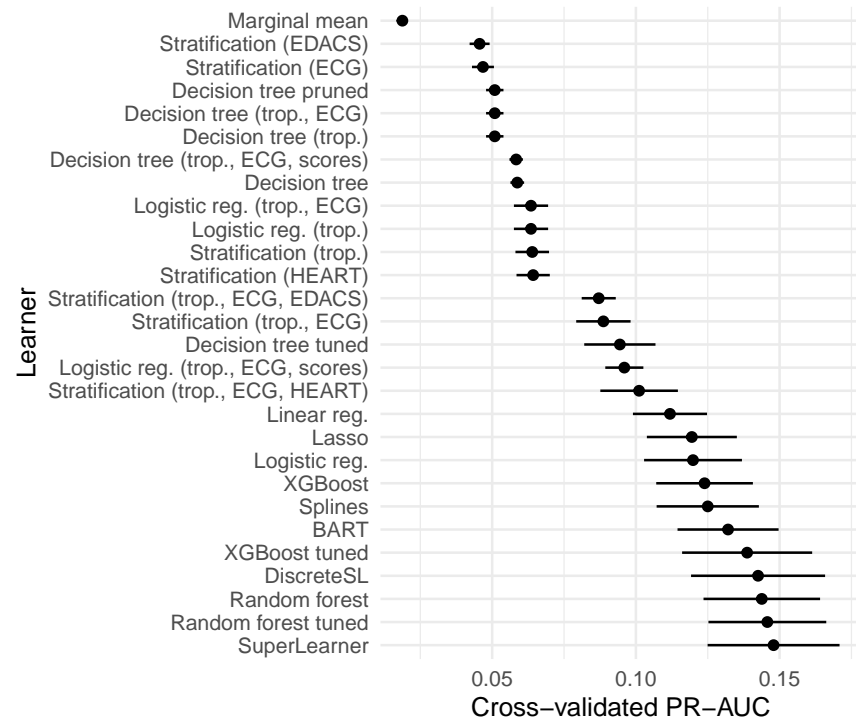
## 3 Results 373

### 3.1 Model performance 374

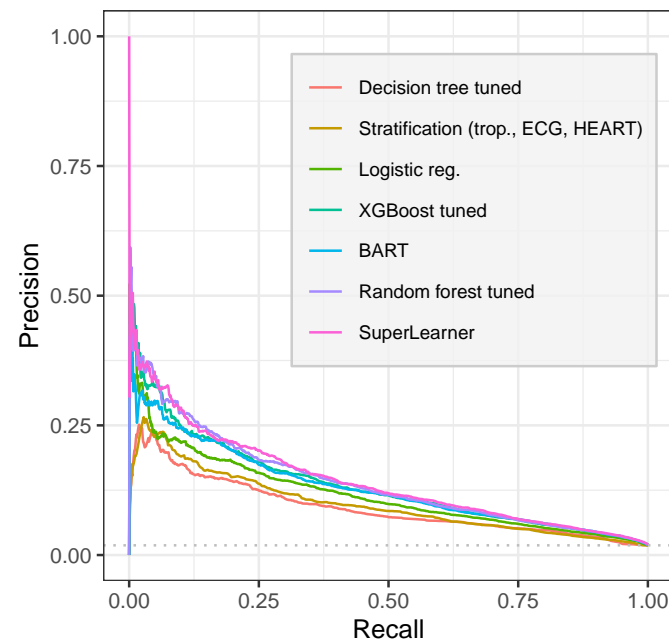
#### 3.1.1 Discrimination 375

Figure 1 displays the estimated precision-recall area under the curve (PR-AUC) PR-AUCs and 95% confidence intervals for each combination of features and estimation algorithm. The MACE mean on the training sample represents the baseline PR-AUC, which was 1.88%. The SuperLearner ensemble achieved the highest estimated PR-AUC (0.148, 95% CI [0.126, 0.170]), followed by the random forest with hyperparameter tuning (0.144, [0.125, 0.164]), the default random forest (0.143, [0.122, 0.165]), and the tuned XGBoost (0.138, [0.116, 0.160]). By comparison the PR-AUC for logistic regression was 0.120 [0.103, 0.137], noticeably lower than the ensemble. Point estimates and confidence intervals are listed in Supplemental Table ?? 376  
377  
378  
379  
380  
381  
382  
383  
384

**Figure 1.** Comparison of cross-validated discriminative performance using PR-AUC metric, with 95% confidence intervals.



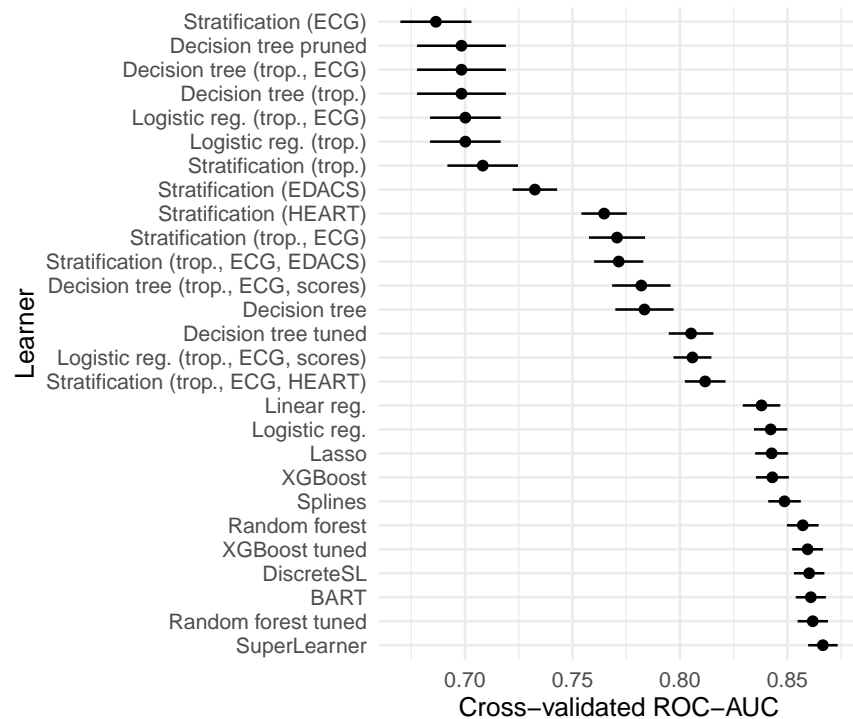
**Figure 2.** Comparison of cross-validated discriminative performance using precision-recall curves for a subset of learners.



For our secondary discrimination metric, cross-validated ROC-AUC was calculated

and is displayed in Figure 3 (LeDell et al. 2015). The SuperLearner ensemble again achieved the highest performance (ROC-AUC = 0.866, 95% CI [0.859, 0.873]), followed by tuned random forest (0.860, [0.853, 0.867]), tuned XGBoost (0.859, [0.852, 0.866]), and BART (0.859, [0.852, 0.866]). The ROC-AUC for logistic regression (0.842, [0.834, 0.850]) was significantly lower than the ensemble ( $p = X$ ). Point estimates and confidence intervals are listed in Supplemental Table ??.

**Figure 3.** Comparison of cross-validated discriminative performance using ROC-AUC metric, with 95% confidence intervals. The simple mean had a standard AUC of 0.5 and is omitted from the plot.



We reviewed the distribution of learner weights in the SuperLearner ensemble to examine which algorithms were used most heavily. The weight distribution of learners that were included at least once (i.e. maximum weight greater than 0) is reported in Table 1. Four algorithms were always incorporated into the ensemble: default random forest (average weight = 0.25), tuned random forest (average weight = 0.25), default XGBoost (average weight = 0.20), and BART (average weight = 0.18). The remaining 3 learners were sometimes incorporated into the ensemble, with average weights ranging from 0.08 to < 0.01 and a maximum individual weight of 0.18.

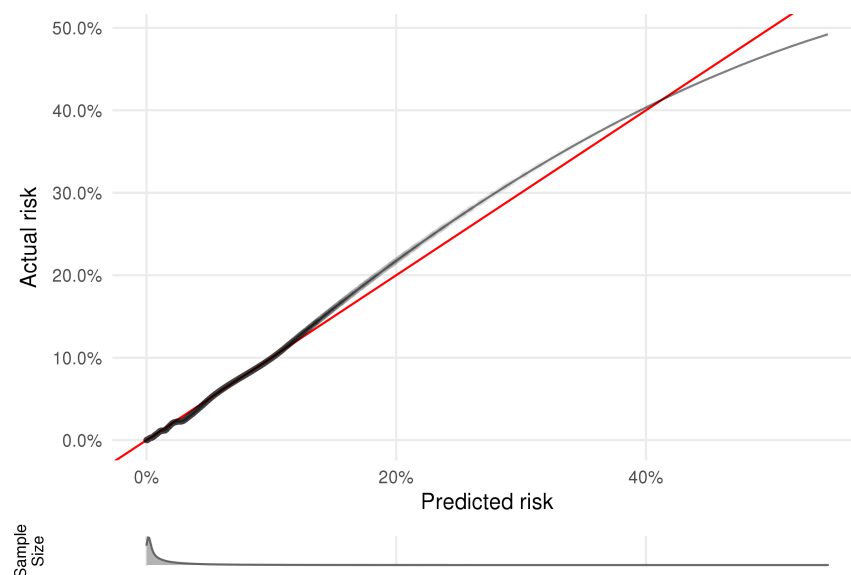
### 3.1.2 Calibration

We visually compared the predicted risk of the SuperLearner ensemble to the lowess-smoothed observed risk in figure 4. The red line is the target calibration in which predicted risk is equal to observed risk. The blue line shows the lowess-smoothed observed risk for each value of the predicted risk.

**Table 1.** Distribution of algorithm weights across ensemble cross-validation replications

#	Learner	Mean	SD	Min	Max
1	Random forest	0.2516	0.0765	0.1815	0.3749
2	Random forest tuned	0.2488	0.0828	0.1449	0.3737
3	XGBoost	0.1959	0.0378	0.1506	0.2378
4	BART	0.1824	0.0535	0.0993	0.2316
5	Splines	0.0748	0.0685	0.0000	0.1764
6	Logistic reg.	0.0429	0.0309	0.0000	0.0720
7	Stratification (trop., ECG, EDACS)	0.0036	0.0080	0.0000	0.0178

**Figure 4.** Calibration plot comparing predicted risk to observed risk

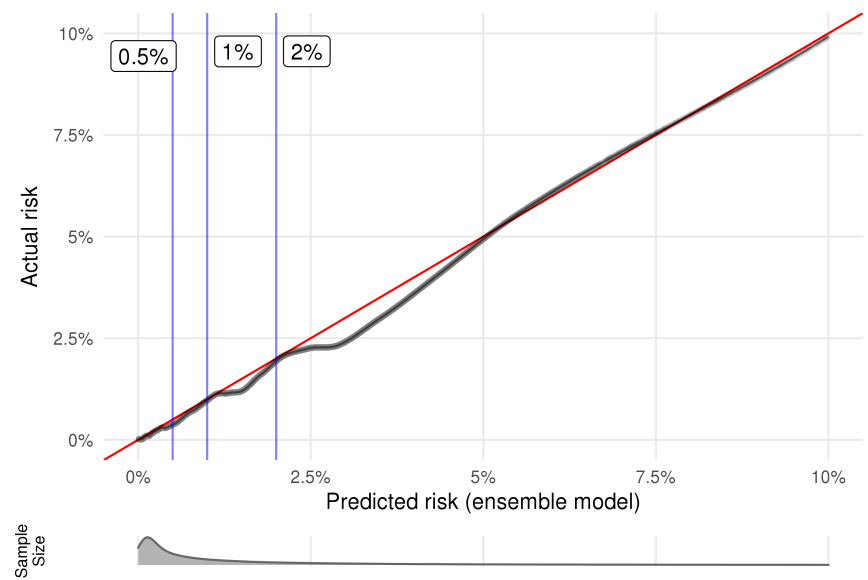


The median predicted risk was 0.64%, with a first quartile of 0.2% and third quartile of 2%. Our primary threshold of scientific interest was 0.5% for possible early discharge. Given those low risk levels, it would be best to “zoom in” our visual calibration review to that region. We show a zoomed calibration plot as Figure 5.

405  
406  
407  
408

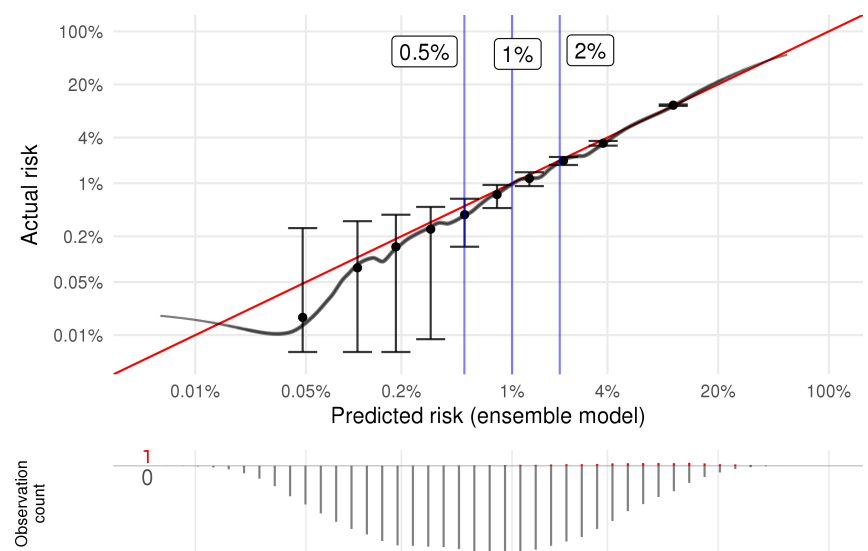


**Figure 5.** Zoomed calibration plot comparing predicted risk to observed risk. Clinical thresholds of 0.5%, 1%, or 2% risk are noted by blue vertical lines.



Finally, we include a exponential-scale calibration plot (Figure 6) with calibration confidence intervals after grouping patients into 10 groups based on predicted risk, consistent with TRIPOD guideline recommendations (Collins et al. 2015). Due to the substantial class imbalance the exponential scaling of axes allows easier comparison across the probability range, although it may be less intuitive due to the shifting of scales. For example, the width of confidence intervals is counterintuitive for visual comparison due to the dynamic scaling, but the amount of information provided is visually consistent throughout the plot.

**Figure 6.** Exponential-scale calibration plot comparing predicted risk to observed risk with grouped 95% confidence intervals. Clinical thresholds of 0.5%, 1%, or 2% risk are noted by blue vertical lines.



**Table 2.** Comparing missing value imputation using GLRM versus median/mode

Variable	Missingness	Error GLRM	Error Median	Percent reduction
HbA1c	59.8	0.088	1.628	94.6
Triglycerides	46.9	0.039	0.528	92.6
HDL	43.5	1.077	14.815	92.7
Total Chol.	42.7	7.961	45.762	82.6
LDL	38.8	4.158	37.291	88.8
GFR	8.8	0.286	11.699	97.6
Trop. 3HV	2.7	0.001	0.008	90.2
ECG	1.8	0.000	0.729	100.0
BMI	1.6	0.645	7.508	91.4
Obese	1.4	0.769	0.640	-20.1
Respiration	0.4	0.020	2.909	99.3
O2 Saturation	0.3	0.020	2.515	99.2
Pulse	0.2	0.680	18.446	96.3
Pulse Peak	0.2	0.703	18.553	96.2
SBP	0.1	0.586	22.843	97.4
Lowest SBP	0.1	0.458	18.673	97.5

As a statistical complement to the visual examination, we also calculated mean absolute error (MAE). MAE is the sample mean of the absolute difference between the smoothed observed risk ( $Risk_O$ ) and the predicted risk ( $Risk_P$ ).

$$\frac{1}{n} \sum_{i=1}^n |Risk_O(i) - Risk_P(i)| \quad (3)$$

We found an MAE of 0.19% with a lowess smoothing span of 0.05 (low smoothing), and an MAE of 0.14% with a smoothing span of 0.20 (moderate smoothing). These statistics indicate that the ensemble risk prediction was typically miscalibrated by about 0.17 percentage points.

### 3.1.3 Missing data imputation

We evaluated the benefit of the more complex GLRM-based imputation by comparing the imputed value to the known value, among variables with missingness. The root mean-squared error metric was calculated for each variable, and for both GLRM and median/mode imputation methods. We could then estimate the percentage improvement in RMSE for the GLRM imputation. Results in Table 2 show a notable improvement in RMSE for every variable, with the exception of the obesity binary variable.

## 3.2 Interpretation

### 3.2.1 Variable importance ranking

As discussed earlier, our objective for the variable importance analysis was to understand which variables were most influential on the prediction of our final model. Providing that ranking could improve the interpretability of the risk prediction, allowing for confirmation that the results are reasonable and possibly yielding additional scientific insights. However, our final model is quite complex: it is a weighted average of multiple versions of random forests, xgboost models, bayesian additive regression trees, etc. In this work we provide rankings for the top two estimation

algorithms: random forest and xgboost. We used the optimal hyperparameter settings from cross-validated analysis.

**Table 3.** Variable importance rankings

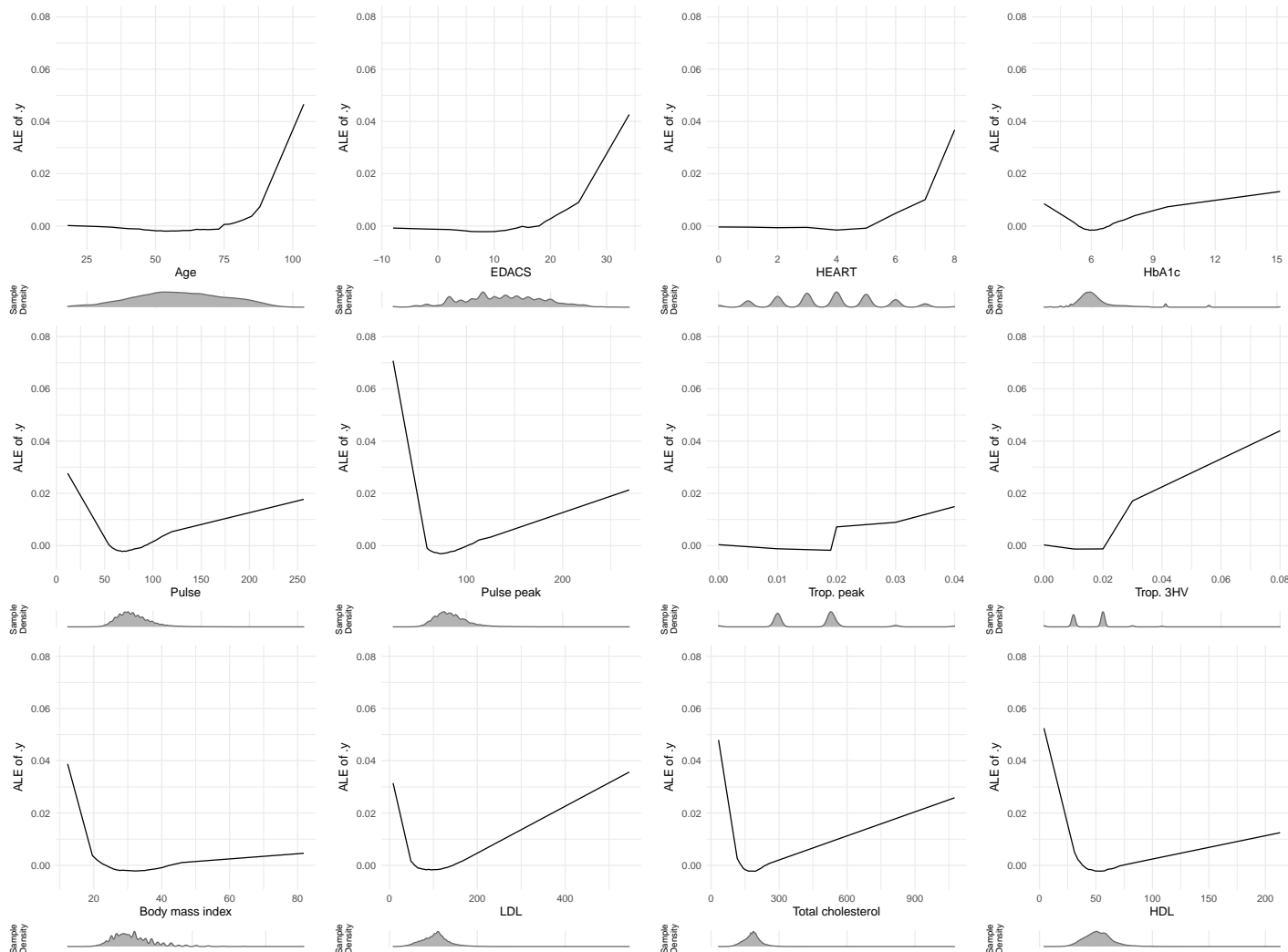
Random Forest importance ranking		XGBoost importance ranking	
Variable	Mean Decrease Accuracy (%)	Variable	Gain
1. Age	0.262	1. Peak troponin	0.3288
2. EDACS	0.188	2. HEART	0.1591
3. HEART	0.117	3. High EDACS	0.0712
4. CAD	0.117	4. EDACS	0.0457
5. Troponin 3HV	0.111	5. High HEART	0.0444
6. LDL	0.104	6. ECG	0.0415
7. Total Cholesterol	0.102	7. Peak pulse	0.0320
8. Missing Triglycerides	0.083	8. Age	0.0282
9. Missing Total Cholesterol	0.080	9. BMI	0.0234
10. Missing LDL	0.076	10. SBP	0.0217
11. Missing HDL	0.073	11. Myocardial infarction	0.0210
12. Pulse	0.071	12. CAD	0.0198
13. Peak pulse	0.070	13. Aortic athero.	0.0175
14. BMI	0.051	14. Troponin 3HV	0.0165
15. Diabetes	0.047	15. HDL	0.0159
16. Hypertension	0.045	16. Respiration	0.0135
17. HDL	0.044	17. O2 saturation	0.0131
18. HbA1c	0.043	18. GFR	0.0130
19. Peak troponin	0.041	19. Exertion	0.0127
20. Triglycerides	0.039	20. Lowest SBP	0.0091

Interestingly we see rather different results between the two algorithms, which supports the hypothesis that an ensemble of multiple algorithms could achieve better performance than selecting a single estimation algorithm. Both algorithms place high emphasis on the EDACS and HEART risk scores, demonstrating the benefit of including those scores along with the underlying predictors. Different versions of the cardiac troponin predictor are emphasized by the two algorithms: random forest focuses on 3-hour troponin whereas xgboost focuses on peak troponin. ECG reading is emphasized by xgboost but not random forest. Both algorithms make use of lipid profile predictors (LDL, HDL) and vital signs (pulse, respiration) that are not included in the existing risk scores. The random forest makes use of certain missingness indicators, which are often indicative of the quality of a patient’s records, while xgboost does not. Also noteworthy is the lack of pain-related characteristics sourced from clinical notes in the top predictors, a difference from prior work highlighting their importance at predicting MACE (Amsterdam et al. 2010).

### 3.2.2 Accumulated local effects

The accumulated local effect method visualizes the conditional relationship of top predictors to the ensemble’s prediction, across their range of values.

**Figure 7.** Accumulated local effect plots of key continuous predictors



## 4 Discussion

459

The next step in model evaluation is to conduct one or more external validations of the discrimination and calibration of the model predictions. This validation might include future retrospective cohorts at the current study location (temporal validation), although preferably cohorts sourced from other regions or EHRs (geographical or institutional validation) (Moons et al. 2012). We hope to collaborate in the future with groups interested in such validations.

460

461

462

463

464

465

In future work we plan to expand the machine learning in several ways. The ensemble weighting could specifically optimize PR-AUC. Incorporating feature selection may benefit the simpler algorithms by removing unhelpful predictors. Feature engineering might be beneficial as well, such as creation of interaction terms or even incorporation of the principal components from the GLRM imputation. Due to computational limitations we were not able to conduct hyperparameter tuning on the BART learner, which likely would provide some performance benefit. We are optimistic that random search or model-based search (e.g. Bayesian optimization) rather than grid search could provide even stronger tuning of algorithm hyperparameters across a higher

466

467

468

469

470

471

472

473

474

number of dimensions. Evaluation of the GLRM imputation could be further contextualized through comparisons to additional imputation methods, especially principal component analysis, k-nearest neighbors, multiple imputation, and variable-specific supervised models (e.g. OLS or random forest). Additional machine learning algorithms could be explored, such as LightGBM, extremely randomized trees, and multivariate adaptive regression splines. The variable importance ranking could be streamlined through a Random Forest-style permutation importance analysis of the SuperLearner ensemble itself, or through a targeted learning method such as vimp (Williamson et al. 2017) or varimpact (Hubbard et al. 2018).

The model might also benefit from a broader sample that includes higher risk patients, which were not included in this study. Calibration might be improved through targeted learning-based adjustment (Brooks et al. 2012). Cross-validated estimation of discrimination performance could be improved through cross-validated targeted maximum likelihood estimation (Benkeser et al. 2019).

## 5 Conclusion

In this work we explored the benefit of complex machine learning algorithms at predicted major adverse cardiac events in patients with chest pain. We found that the ML algorithms were able to achieve improved discrimination compared to simpler baselines such as logistic regression, decision trees, or stratification on individual predictors. Combining multiple algorithms into an ensemble estimator yielded the best performance, and rather than select optimal hyperparameters we created an ensemble of algorithms across different hyperparameters. We demonstrated the surprising effectiveness of generalized low-rank models for imputation of missingness in EHR-sourced patient data. Finally, we provided interpretation of how the ensemble's prediction is generated through two methods: ranking the predictors by their contribution to predictive performance, and visualizing the dose-response effect of continuous predictors with accumulated local effect plots.

The cleaning and analysis code for this project has been translated to use a public dataset and is available online at <https://github.com/ck37/Predictive-Modeling-in-R>. Functions to calculate PR-AUC, ROC-AUC, index of prediction accuracy (IPA), and Brier scores for cross-validated SuperLearner ensembles are provided in the open source R package ck37r (C. J. Kennedy 2020).

## Funding

This work was supported by a Kaiser Permanente Division of Research Delivery Science Research Grant.

## Acknowledgments

We thank Dustin Ballard, Gabriel Escobar, Alan Ho, Oleg Sofrygin, and Jodi McCloskey for helpful comments. We thank Adina Rauchwerger and Laura Simon for project management.

## References

- Agor, Joseph, Osman Y Özaltın, Julie S Ivy, Muge Capan, Ryan Arnold, and Santiago Romero (2019). “The value of missing information in severity of illness score development”. In: *Journal of biomedical informatics* 97, p. 103255.
- Amsterdam, Ezra A, J Douglas Kirk, David A Bluemke, Deborah Diercks, Michael E Farkouh, J Lee Garvey, Michael C Kontos, James McCord, Todd D Miller, Anthony Morise, et al. (2010). “Testing of low-risk patients presenting to the emergency department with chest pain: a scientific statement from the American Heart Association”. In: *Circulation* 122.17, pp. 1756–1776.
- Apley, Daniel W and Jingyu Zhu (2019). “Visualizing the effects of predictor variables in black box supervised learning models”. In: *arXiv preprint arXiv:1612.08468*.
- Austin, Peter C (2007). “A comparison of regression trees, logistic regression, generalized additive models, and multivariate adaptive regression splines for predicting AMI mortality”. In: *Statistics in medicine* 26.15, pp. 2937–2957.
- Austin, Peter C and Jack V Tu (2004). “Bootstrap methods for developing predictive models”. In: *The American Statistician* 58.2, pp. 131–137.
- Benkeser, David, Maya Petersen, and Mark J van der Laan (2019). “Improved small-sample estimation of nonlinear cross-validated prediction metrics”. In: *Journal of the American Statistical Association*, pp. 1–16.
- Bonaca, Marc, Benjamin Scirica, Marc Sabatine, Anthony Dalby, Jindrich Spinar, Sabina A Murphy, Peter Jarolim, Eugene Braunwald, and David A Morrow (2010). “Prospective evaluation of the prognostic implications of improved assay performance with a sensitive assay for cardiac troponin I”. In: *Journal of the American College of Cardiology* 55.19, pp. 2118–2124.
- Breiman, Leo (1996). “Stacked regressions”. In: *Machine learning* 24.1, pp. 49–64.
- (2001). “Random forests”. In: *Machine learning* 45.1, pp. 5–32.
- Brier, Glenn W (1950). “Verification of forecasts expressed in terms of probability”. In: *Monthly weather review* 78.1, pp. 1–3.
- Brooks, Jordan, Mark J van der Laan, and Alan S Go (2012). “Targeted maximum likelihood estimation for prediction calibration”. In: *The international journal of biostatistics* 8.1.
- Chen, Tianqi and Carlos Guestrin (2016). “Xgboost: A scalable tree boosting system”. In: *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794.
- Chipman, Hugh A, Edward I George, Robert E McCulloch, et al. (2010). “BART: Bayesian additive regression trees”. In: *The Annals of Applied Statistics* 4.1, pp. 266–298.
- Christodoulou, Evangelia, MA Jie, Gary S Collins, Ewout W Steyerberg, Jan Y Verbakel, Ben van Calster, et al. (2019). “A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models”. In: *Journal of clinical epidemiology*.
- Collins, Gary S, Johannes B Reitsma, Douglas G Altman, and Karel GM Moons (2015). “Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD statement”. In: *British Journal of Surgery* 102.3, pp. 148–158.
- Cook, Nancy R (2007). “Use and misuse of the receiver operating characteristic curve in risk prediction”. In: *Circulation* 115.7, pp. 928–935.
- De Lemos, James A, Mark H Drazner, Torbjorn Omland, Colby R Ayers, Amit Khera, Anand Rohatgi, Ibrahim Hashim, Jarett D Berry, Sandeep R Das, David A Morrow, et al. (2010). “Association of troponin T detected with a highly sensitive assay and cardiac structure and mortality risk in the general population”. In: *Jama* 304.22, pp. 2503–2512.



- Friedman, Jerome H (2001). “Greedy function approximation: a gradient boosting machine”. In: *Annals of statistics*, pp. 1189–1232.
- Goldstein, Benjamin A, Ann Marie Navar, and Rickey E. Carter (2016). “Moving beyond regression techniques in cardiovascular risk prediction: applying machine learning to address analytic challenges”. In: *European Heart Journal* 38.23, pp. 1805–1814. ISSN: 0195-668X. DOI: [10.1093/eurheartj/ehw302](https://doi.org/10.1093/eurheartj/ehw302).
- Goldstein, Benjamin A, Ann Marie Navar, Michael J Pencina, and John P A Ioannidis (2016). “Opportunities and challenges in developing risk prediction models with electronic health records data: a systematic review”. In: *Journal of the American Medical Informatics Association* 24.1, pp. 198–208. ISSN: 1067-5027. DOI: [10.1093/jamia/ocw042](https://doi.org/10.1093/jamia/ocw042).
- Greenslade, Jaimi H, Edward W Carlton, Christopher Van Hise, Elizabeth Cho, Tracey Hawkins, William A Parsonage, Jillian Tate, Jacobus Ungerer, and Louise Cullen (2018). “Diagnostic accuracy of a new high-sensitivity troponin I assay and five accelerated diagnostic pathways for ruling out acute myocardial infarction and acute coronary syndrome”. In: *Annals of emergency medicine* 71.4, pp. 439–451.
- Hastie, Trevor J and Robert J Tibshirani (1990). *Generalized additive models*. Vol. 43. CRC press.
- Hastie, Trevor, Robert Tibshirani, and Ryan J Tibshirani (2017). “Extended comparisons of best subset selection, forward stepwise selection, and the lasso”. In: *arXiv preprint arXiv:1707.08692*.
- He, Jianxing, Sally L Baxter, Jie Xu, Jiming Xu, Xingtao Zhou, and Kang Zhang (2019). “The practical implementation of artificial intelligence technologies in medicine”. In: *Nature medicine* 25.1, p. 30.
- Hilden, Jørgen and Thomas A Gerds (2014). “A note on the evaluation of novel biomarkers: do not rely on integrated discrimination improvement and net reclassification index”. In: *Statistics in medicine* 33.19, pp. 3405–3414.
- Hubbard, Alan E, Chris J Kennedy, and Mark J van der Laan (2018). “Data-Adaptive Target Parameters”. In: *Targeted Learning in Data Science*. Springer, pp. 125–142.
- Janssens, A Cecile J W and Forike K Martens (2020). “Reflection on modern methods: Revisiting the area under the ROC Curve”. In: *International Journal of Epidemiology*. dyz274. ISSN: 0300-5771. DOI: [10.1093/ije/dyz274](https://doi.org/10.1093/ije/dyz274).
- Johnson, Kipp W, Jessica Torres Soto, Benjamin S Glicksberg, Khader Shameer, Riccardo Miotto, Mohsin Ali, Euan Ashley, and Joel T Dudley (2018). “Artificial intelligence in cardiology”. In: *Journal of the American College of Cardiology* 71.23, pp. 2668–2679.
- Kattan, Michael W and Thomas A Gerds (2018). “The index of prediction accuracy: an intuitive measure useful for evaluating risk prediction models”. In: *Diagnostic and prognostic research* 2.1, p. 7.
- Kennedy, Chris J (2020). *ck37r: Chris Kennedy’s R toolkit*. URL: <https://github.com/ck37/ck37r>.
- Kennedy, Edward H, Wyndy L Wiitala, Rodney A Hayward, and Jeremy B Sussman (2013). “Improved cardiovascular risk prediction using nonparametric regression and electronic health record data”. In: *Medical care* 51.3, p. 251.
- Kerr, Kathleen F, Zheyu Wang, Holly Janes, Robyn L McClelland, Bruce M Psaty, and Margaret S Pepe (2014). “Net reclassification indices for evaluating risk-prediction instruments: a critical review”. In: *Epidemiology (Cambridge, Mass.)* 25.1, p. 114.
- Kramer, Andrew A and Jack E Zimmerman (2007). “Assessing the calibration of mortality benchmarks in critical care: The Hosmer-Lemeshow test revisited”. In: *Critical care medicine* 35.9, pp. 2052–2056.

- LeDell, Erin, Maya Petersen, and Mark van der Laan (2015). “Computationally efficient confidence intervals for cross-validated area under the ROC curve estimates”. In: *Electronic journal of statistics* 9.1, p. 1583.
- Leening, Maarten JG, Moniek M Vedder, Jacqueline CM Witteman, Michael J Pencina, and Ewout W Steyerberg (2014). “Net reclassification improvement: computation, interpretation, and controversies: a literature review and clinician’s guide”. In: *Annals of internal medicine* 160.2, pp. 122–131.
- Lichtenstein, Sarah, Baruch Fischhoff, and Lawrence D Phillips (1981). *Calibration of probabilities: The state of the art to 1980*. Tech. rep. Decision Research. Eugene, OR.
- Mark, Dustin G, Jie Huang, Uli Chettipally, Mamata V Kene, Megan L Anderson, Erik P Hess, Dustin W Ballard, David R Vinson, and Mary E Reed (2018). “Performance of coronary risk scores among patients with chest pain in the emergency department”. In: *Journal of the American College of Cardiology* 71.6, pp. 606–616.
- Molnar, Christoph (2020). *Interpretable Machine Learning*. Lulu. com.
- Moons, Karel GM, Andre Pascal Kengne, Diederick E Grobbee, Patrick Royston, Yvonne Vergouwe, Douglas G Altman, and Mark Woodward (2012). “Risk prediction models: II. External validation, model updating, and impact assessment”. In: *Heart* 98.9, pp. 691–698.
- Murphy, Allan H and Robert L Winkler (1977). “Reliability of subjective probability forecasts of precipitation and temperature”. In: *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 26.1, pp. 41–47.
- Pepe, Margaret S, Jing Fan, Ziding Feng, Thomas Gerds, and Jorgen Hilden (2015). “The net reclassification index (NRI): a misleading measure of prediction improvement even with independent test data sets”. In: *Statistics in biosciences* 7.2, pp. 282–295.
- Probst, Philipp, Anne-Laure Boulesteix, and Bernd Bischl (2019). “Tunability: Importance of Hyperparameters of Machine Learning Algorithms.” In: *Journal of Machine Learning Research* 20.53, pp. 1–32.
- Probst, Philipp, Marvin N Wright, and Anne-Laure Boulesteix (2019). “Hyperparameters and tuning strategies for random forest”. In: *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 9.3, e1301.
- Rozenholc, Yves, Thoralf Mildenerger, and Ursula Gather (2010). “Combining regular and irregular histograms by penalized likelihood”. In: *Computational Statistics & Data Analysis* 54.12, pp. 3313–3323.
- Rui, P, K Kang, and JJ. Ashman (2016). *National Hospital Ambulatory Medical Care Survey: 2016 emergency department summary tables*. URL: [https://www.cdc.gov/nchs/data/ahcd/nhamcs\\_emergency/2016\\_ed\\_web\\_tables.pdf](https://www.cdc.gov/nchs/data/ahcd/nhamcs_emergency/2016_ed_web_tables.pdf).
- Saito, Takaya and Marc Rehmsmeier (2015). “The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets”. In: *PloS one* 10.3.
- Sanders, Frederick (1963). “On subjective probability forecasting”. In: *Journal of Applied Meteorology* 2.2, pp. 191–201.
- Schuler, Alejandro, Vincent Liu, Joe Wan, Alison Callahan, Madeleine Udell, David E Stark, and Nigam H Shah (2016). “Discovering patient phenotypes using generalized low rank models”. In: *Biocomputing 2016: Proceedings of the Pacific Symposium*. World Scientific, pp. 144–155.
- Segal, Mark and Yuanyuan Xiao (2011). “Multivariate random forests”. In: *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 1.1, pp. 80–87.
- Senn, Stephen (2005). “Dichotomania: an obsessive compulsive disorder that is badly affecting the quality of analysis of pharmaceutical trials”. In: *Proceedings of the International Statistical Institute, 55th Session, Sydney*.

- Steyerberg, Ewout W (2009). *Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating*. Springer Science & Business Media.
- Steyerberg, Ewout W, Frank E Harrell Jr, Gerard JJM Borsboom, MJC Eijkemans, Yvonne Vergouwe, and J Dik F Habbema (2001). “Internal validation of predictive models: efficiency of some procedures for logistic regression analysis”. In: *Journal of clinical epidemiology* 54.8, pp. 774–781.
- Steyerberg, Ewout W, Andrew J Vickers, Nancy R Cook, Thomas Gerds, Mithat Gonen, Nancy Obuchowski, Michael J Pencina, and Michael W Kattan (2010). “Assessing the performance of prediction models: a framework for some traditional and novel measures”. In: *Epidemiology (Cambridge, Mass.)* 21.1, p. 128.
- Than, Martin, Dylan Flaws, Sharon Sanders, Jenny Doust, Paul Glasziou, Jeffery Kline, Sally Aldous, Richard Troughton, Christopher Reid, William A Parsonage, et al. (2014). “Development and validation of the Emergency Department Assessment of Chest Pain Score and 2 h accelerated diagnostic protocol”. In: *Emergency Medicine Australasia* 26.1, pp. 34–44.
- Than, Martin, Mel Herbert, Dylan Flaws, Louise Cullen, Erik Hess, Judd E Hollander, Deborah Diercks, Michael W Ardagh, Jeffery A Kline, Zea Munro, et al. (2013). “What is an acceptable risk of major adverse cardiac event in chest pain patients soon after discharge from the Emergency Department?: a clinical survey”. In: *International journal of cardiology* 166.3, pp. 752–754.
- Tibshirani, Robert (1996). “Regression shrinkage and selection via the lasso”. In: *Journal of the Royal Statistical Society: Series B (Methodological)* 58.1, pp. 267–288.
- Udell, Madeleine, Corinne Horn, Reza Zadeh, Stephen Boyd, et al. (2016). “Generalized low rank models”. In: *Foundations and Trends® in Machine Learning* 9.1, pp. 1–118.
- Van Calster, Ben, David J McLernon, Maarten van Smeden, Laure Wynants, and Ewout W Steyerberg (2019). “Calibration: the Achilles heel of predictive analytics”. In: *BMC medicine* 17.1, pp. 1–7.
- van der Laan, Mark J, Eric C Polley, and Alan E Hubbard (2007). “Super learner”. In: *Statistical applications in genetics and molecular biology* 6.1.
- Wang, R, J Sedransk, and JH Jinn (1992). “Secondary data analysis when there are missing observations”. In: *Journal of the American Statistical Association* 87.420, pp. 952–961.
- Williamson, Brian D, Peter B Gilbert, Noah Simon, and Marco Carone (2017). “Nonparametric variable importance assessment using machine learning techniques”. In:
- Wolpert, David H (1992). “Stacked generalization”. In: *Neural networks* 5.2, pp. 241–259.
- Wood, Simon N (2003). “Thin plate regression splines”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 65.1, pp. 95–114.
- Yates, J Frank (1982). “External correspondence: Decompositions of the mean probability score”. In: *Organizational Behavior and Human Performance* 30.1, pp. 132–156.