

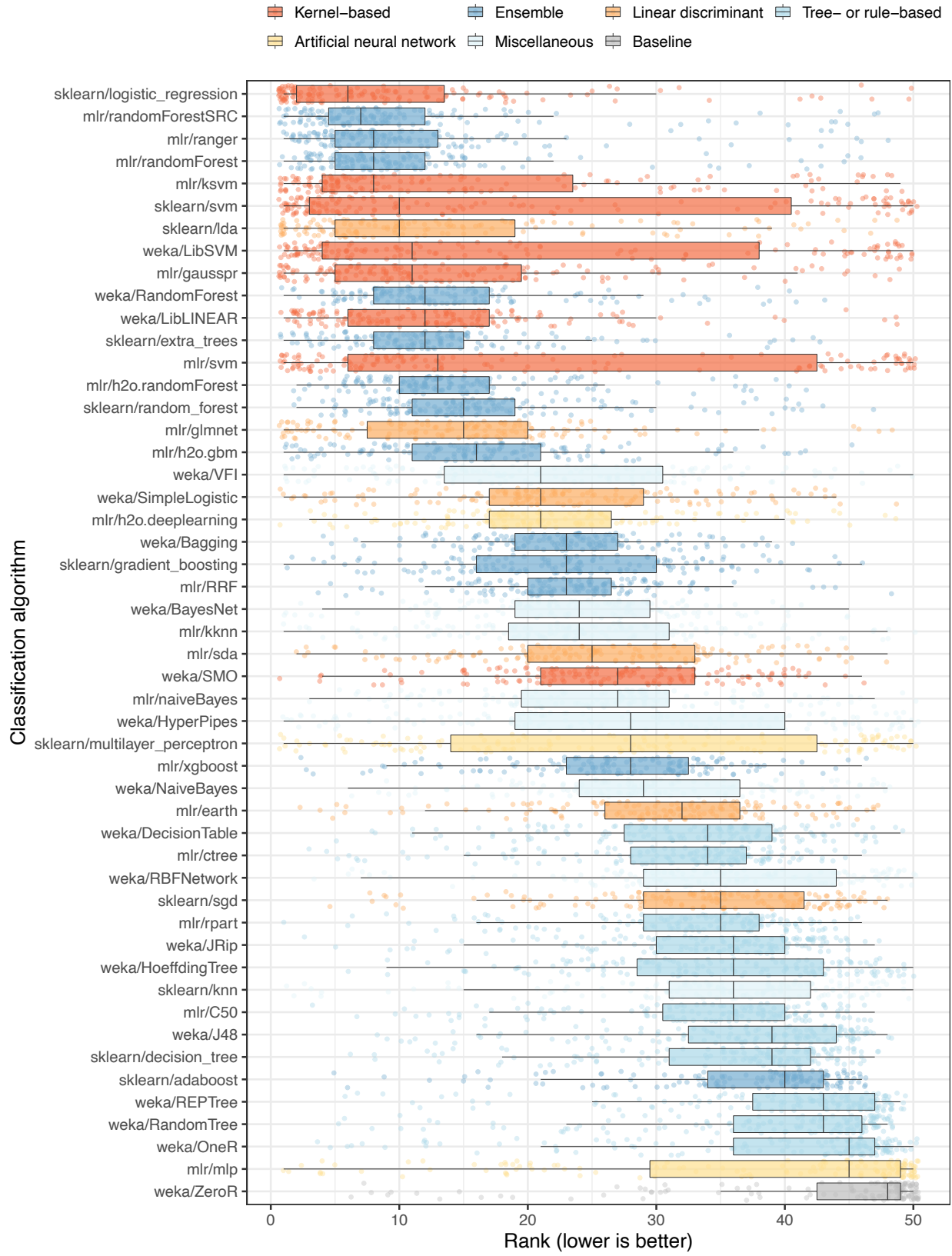
1 **Benchmarking 50 classification algorithms on 50 gene-**
2 **expression datasets**

3 Stephen R. Piccolo^{1,*}, Avery Mecham¹, Nathan P. Golightly¹, Jérémie L. Johnson¹, Dustin B. Miller¹

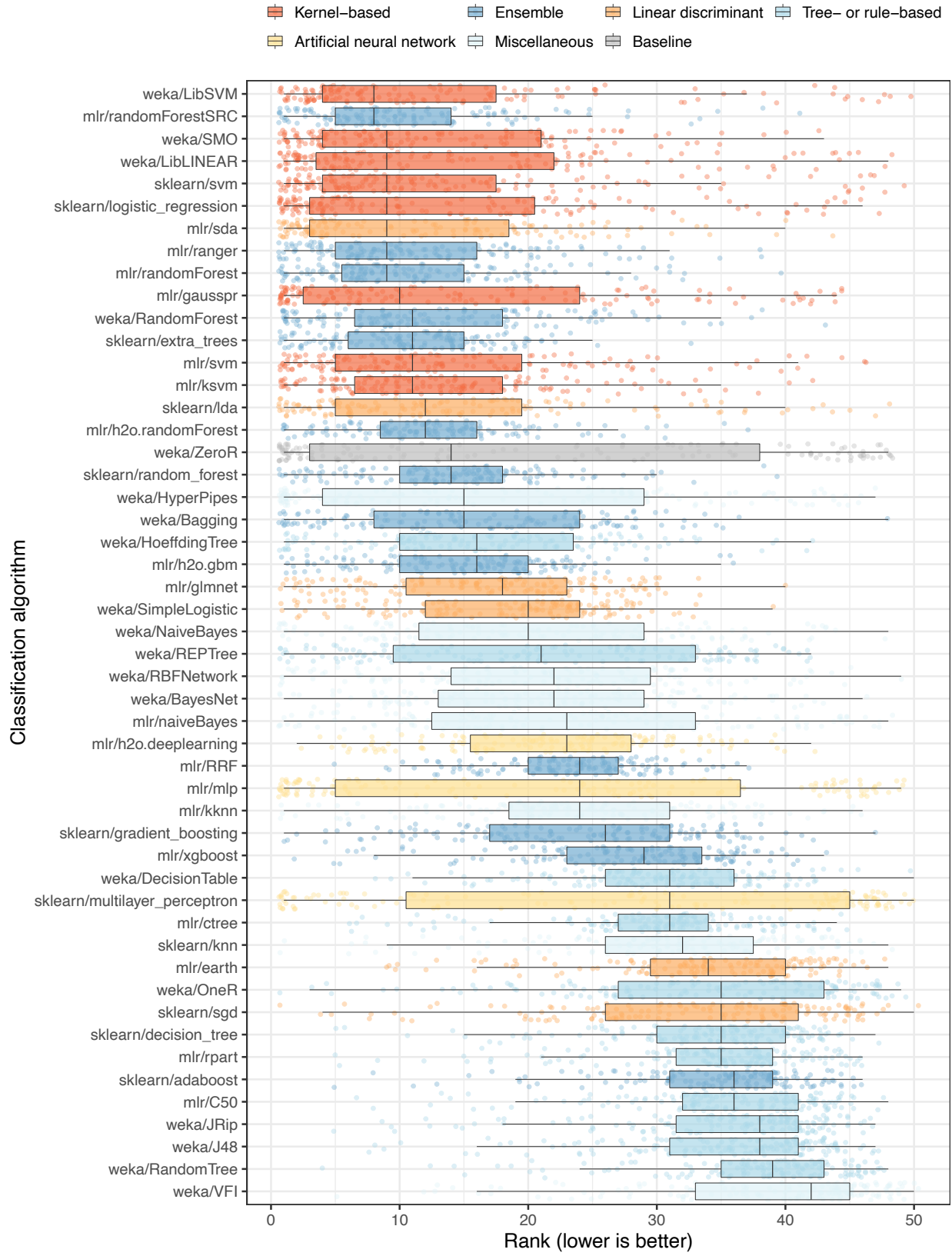
4 1 - Department of Biology, Brigham Young University, Provo, UT, USA

5 * - Please address correspondence to S.R.P. at stephen_piccolo@byu.edu.

6 **Supplementary Figures**



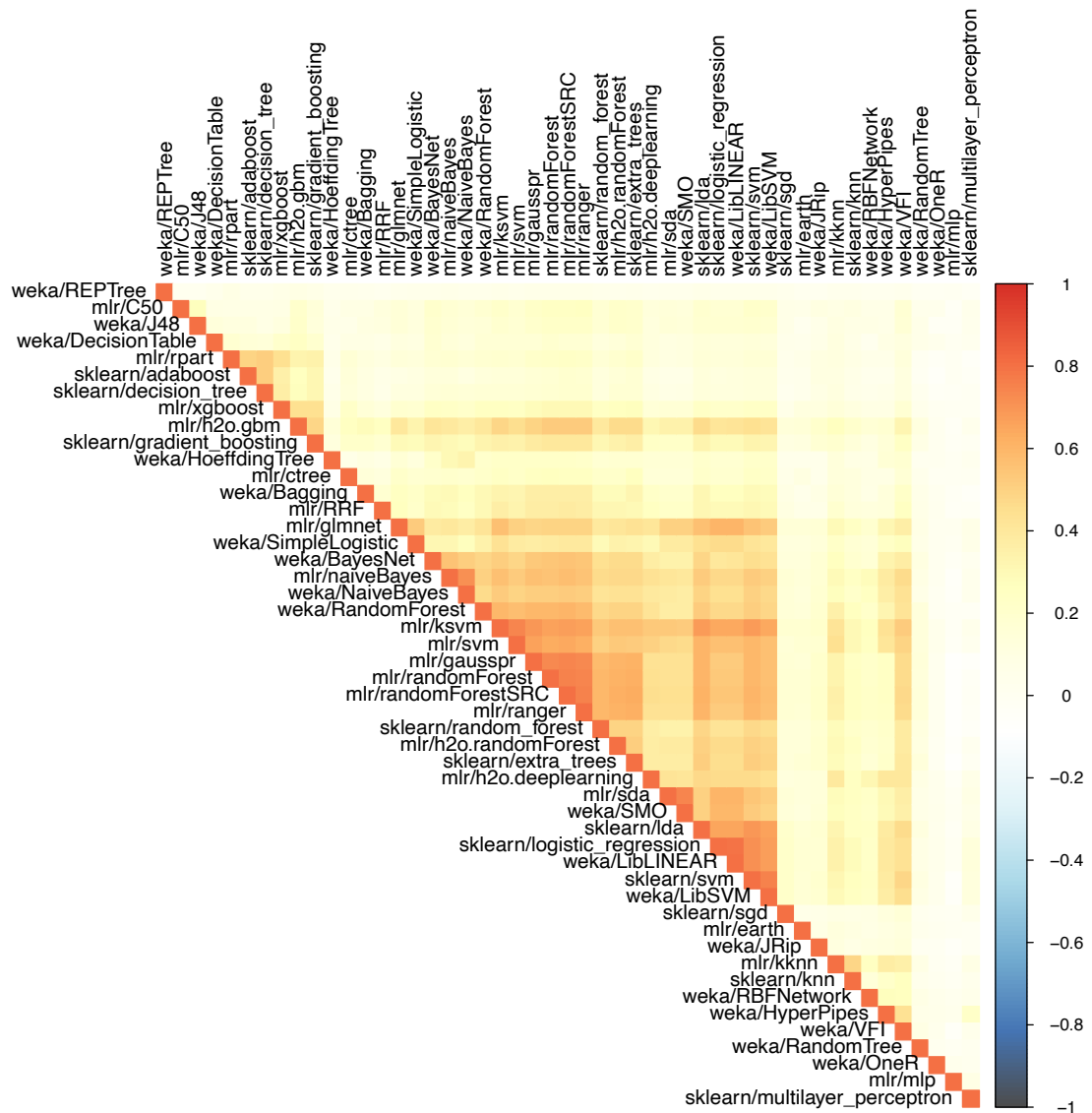
8 **Figure S1: Relative performance of classification algorithms using gene-expression predictors and**
9 **area under the receiver operating characteristic curve as the metric.** We predicted patient states using
10 gene-expression predictors only (Analysis 1). For each combination of dataset, class variable, and
11 classification algorithm, we calculated the arithmetic mean of area under the receiver operating
12 characteristic curve (AUROC) values across 50 iterations of Monte Carlo cross-validation. Next we
13 sorted the algorithms based on the average rank across all dataset/class combinations. Each data point that
14 overlays the box plots represents a particular dataset/class combination. The top 15 performers (relatively
15 low ranks) were algorithms that use linear decision boundaries, kernel functions, and/or ensembles of
16 decision trees.



18 **Figure S2: Relative performance of classification algorithms using gene-expression predictors and**
19 **classification accuracy as the metric.** We predicted patient states using gene-expression predictors only
20 (Analysis 1). For each combination of dataset, class variable, and classification algorithm, we calculated
21 the arithmetic mean of classification accuracy across 50 iterations of Monte Carlo cross-validation. Next
22 we sorted the algorithms based on the average rank across all dataset/class combinations. Each data point
23 that overlays the box plots represents a particular dataset/class combination.



25 **Figure S3: Relative performance of classification algorithms using gene-expression predictors and**
26 **Matthews Correlation Coefficient as the metric.** We predicted patient states using gene-expression
27 predictors only (Analysis 1). For each combination of dataset, class variable, and classification algorithm,
28 we calculated the arithmetic mean of the Matthews Correlation Coefficient across 50 iterations of Monte
29 Carlo cross-validation. Next we sorted the algorithms based on the average rank across all dataset/class
30 combinations. Each data point that overlays the box plots represents a particular dataset/class
31 combination.



32

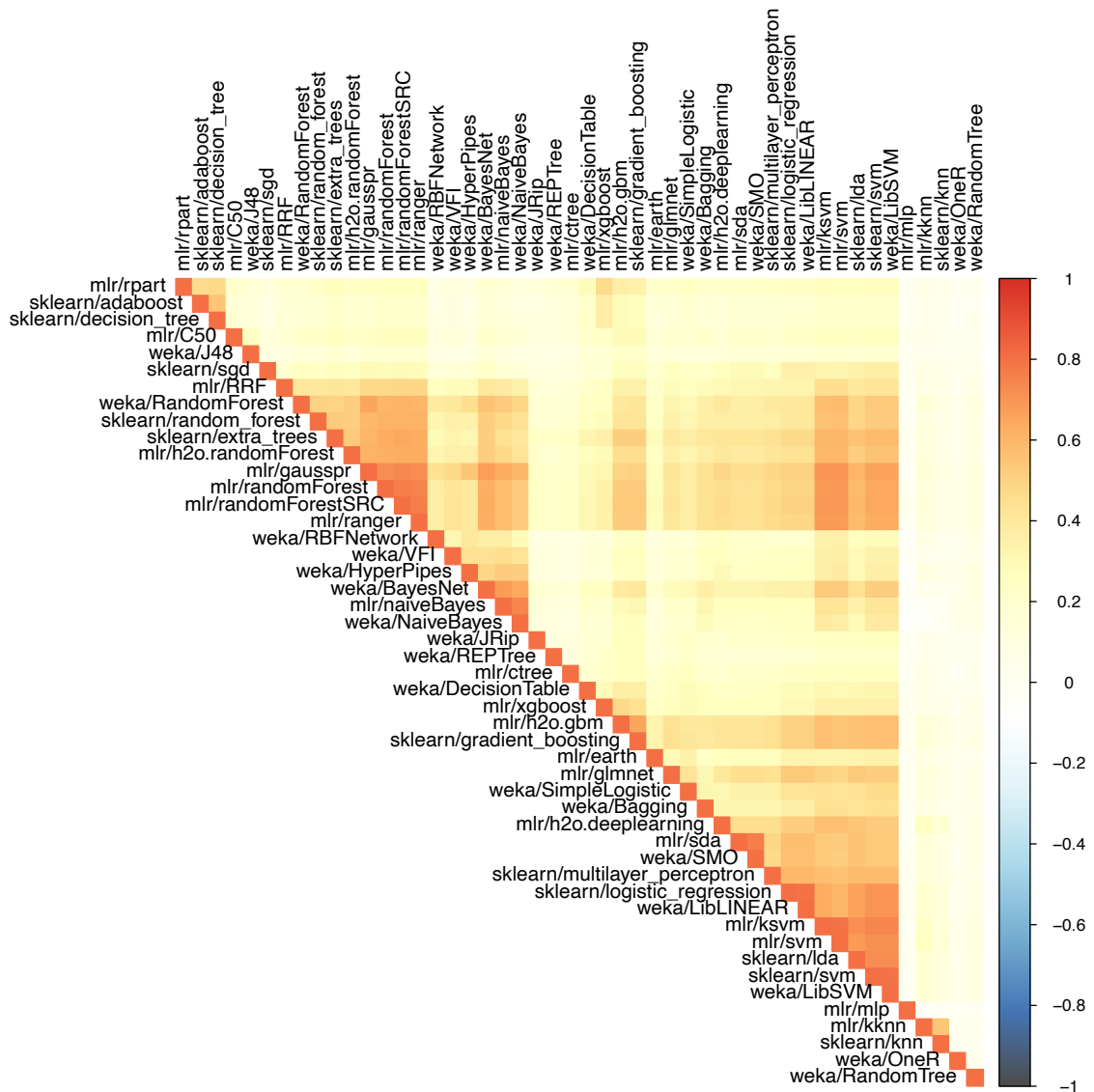
33 **Figure S4: Pairwise correlations of sample-level, probabilistic predictions between classification**

34 **algorithms for dataset GSE10320.** We used each classification algorithm to make probabilistic

35 predictions of relapse in Wilms tumor patients (GSE10320). Based on these predictions, we calculated the

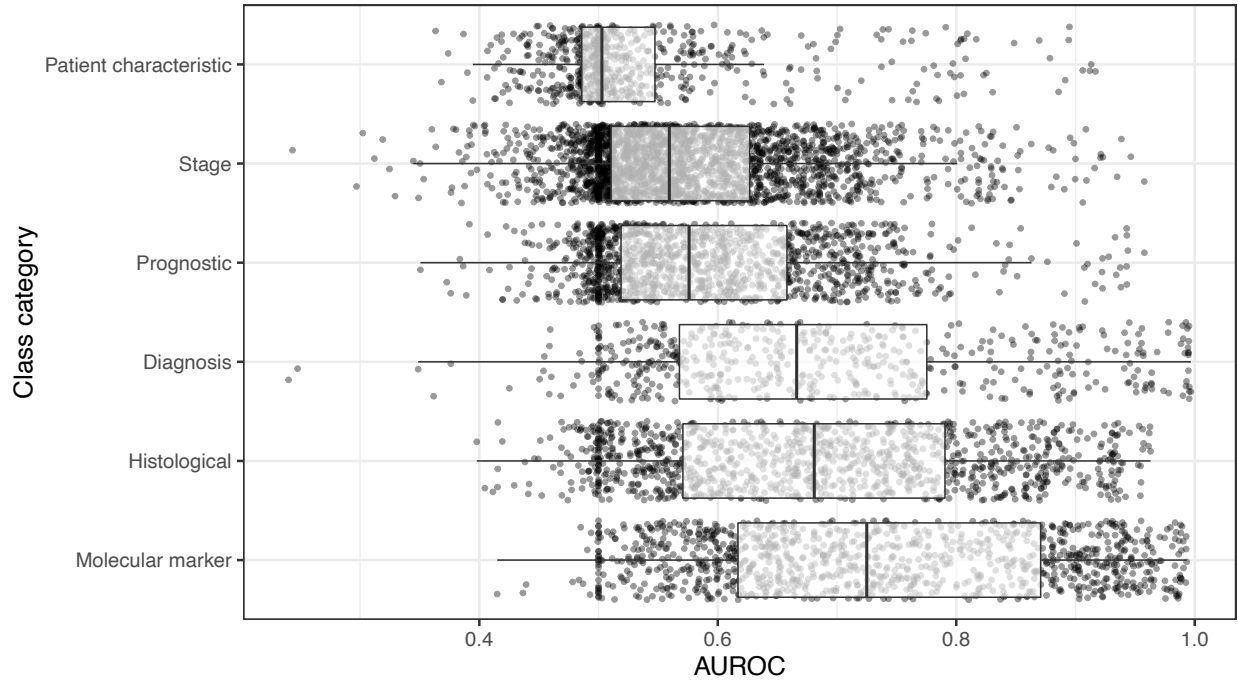
36 Spearman correlation coefficient for each pair of algorithms. These coefficients, averaged across Monte

37 Carlo cross-validation iterations, are illustrated as a correlation plot, clustered based on similarity.



38

39 **Figure S5: Pairwise correlations of sample-level, probabilistic predictions between classification**
 40 **algorithms for dataset GSE46691.** We used each classification algorithm to make probabilistic
 41 predictions of early metastasis following radical prostatectomy (GSE46691). Based on these predictions,
 42 we calculated the Spearman correlation coefficient for each pair of algorithms. These coefficients,
 43 averaged across Monte Carlo cross-validation iterations, are illustrated as a correlation plot, clustered
 44 based on similarity.



45

46 **Figure S6: Dataset performance by class category when using gene-expression predictors.** For each

47 class variable across all datasets, we assigned a category representing the type of patient state being

48 predicted. For Analysis 1, we show the predictive performance for each combination of dataset, class

49 variable, and classification algorithm in each class category. We use area under the receiver operating

50 characteristic curve (AUROC) as the metric. The dashed, red line indicates the performance expected by

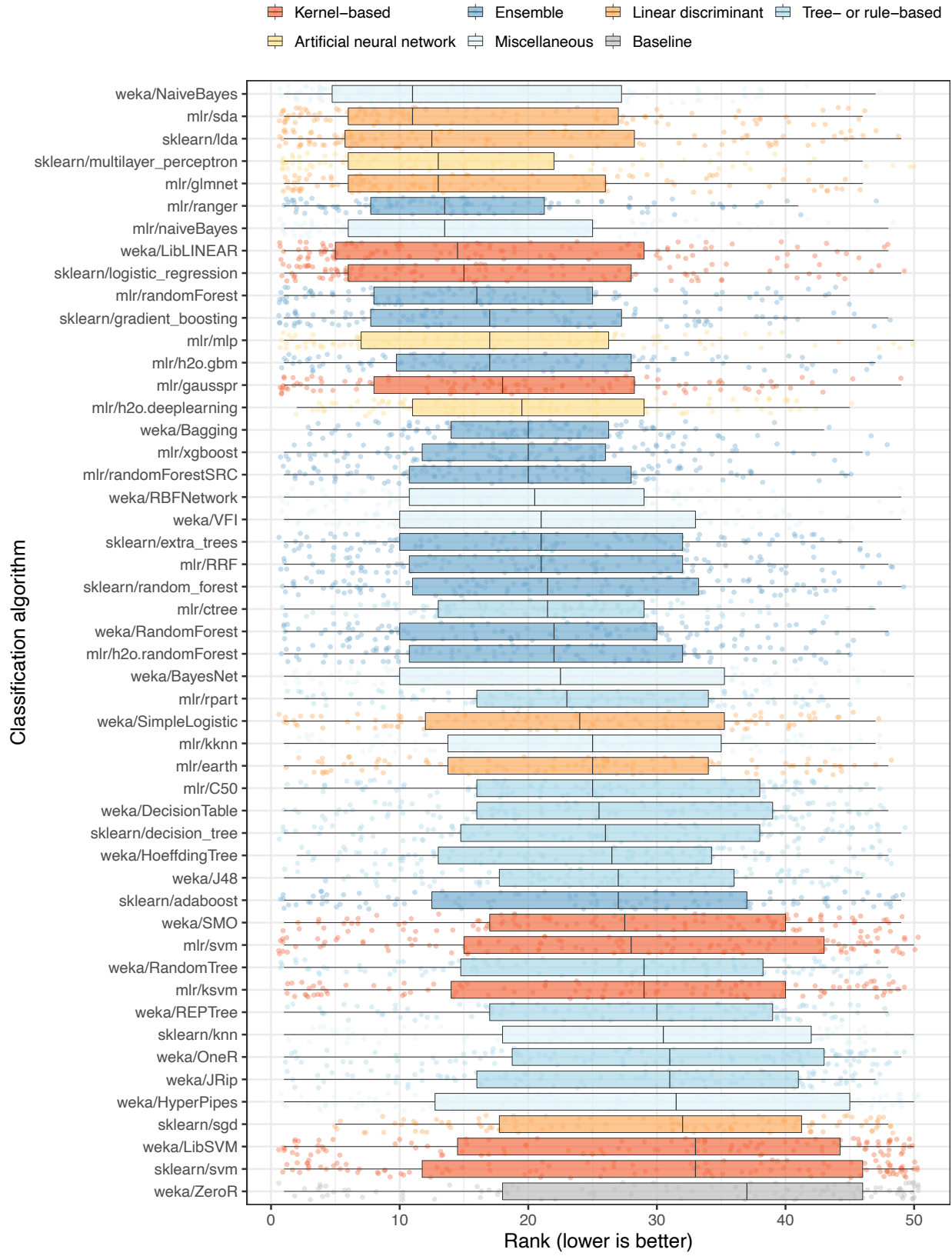
51 random chance. The top-performing category was “Molecular Marker,” which includes class variables

52 associated with mutation status, immunohistochemistry markers of protein expression, presence or

53 absence of chromosomal aberrations, etc. The lowest-performing category was “Patient Characteristic,”

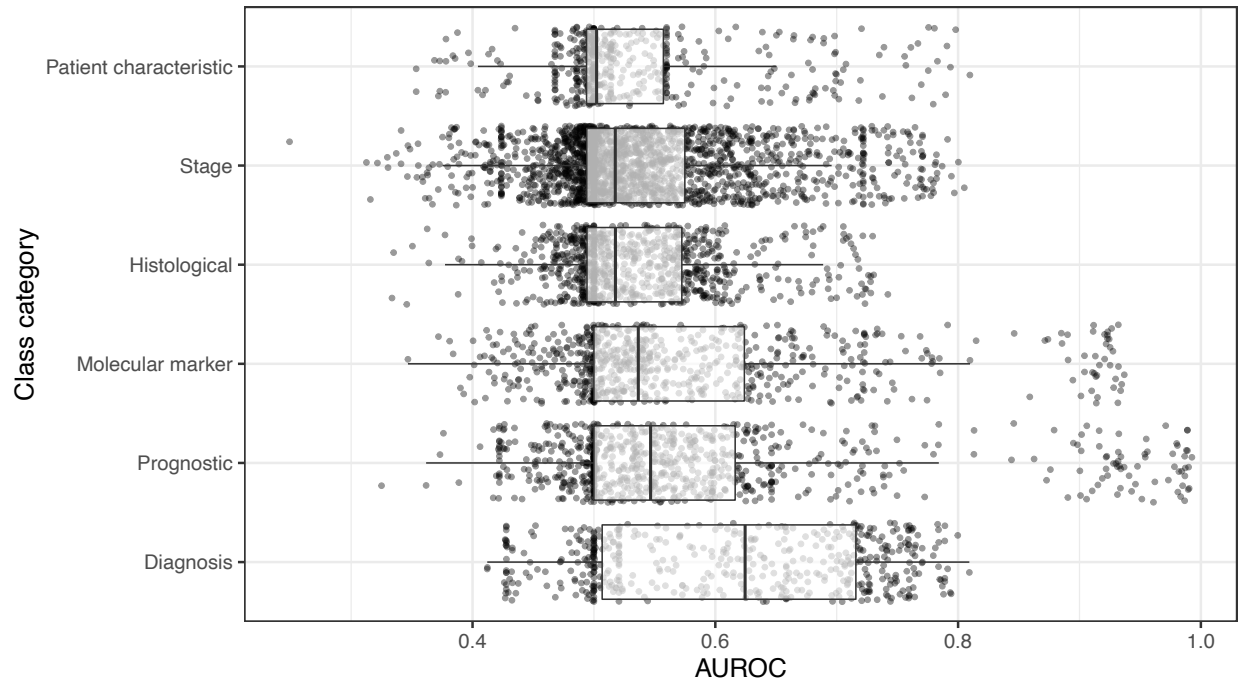
54 which includes variables that indicate whether patients had a family history of cancer, had been diagnosed

55 with multiple tumors, patient performance status, etc.



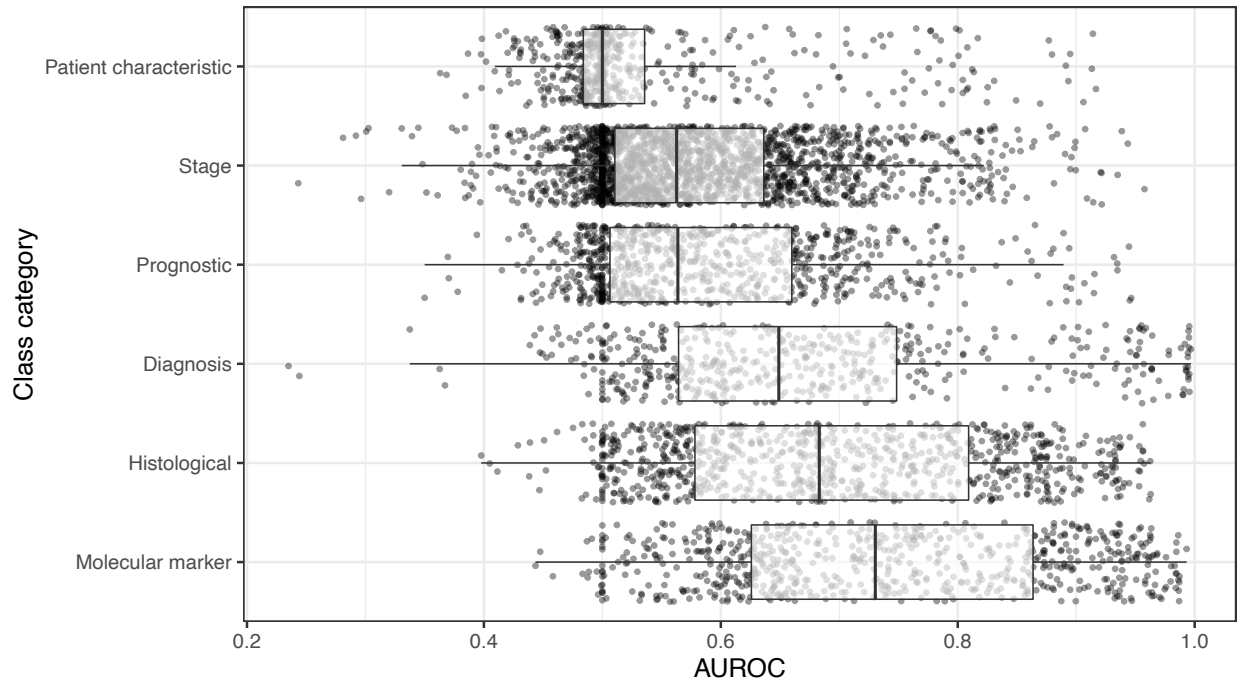
57 **Figure S7: Relative performance of classification algorithms using clinical predictors and area**
58 **under the receiver operating characteristic curve as the metric.** We predicted patient states using
59 clinical predictors only (Analysis 2). For each combination of dataset, class variable, and classification
60 algorithm, we calculated the arithmetic mean of area under the receiver operating characteristic curve
61 (AUROC) values across 50 iterations of Monte Carlo cross-validation. Next we sorted the algorithms
62 based on the average rank across all dataset/class combinations. Each data point that overlays the box
63 plots represents a particular dataset/class combination (some datasets did not have clinical predictors).
64 The top-performing algorithms (relatively low ranks) were similar overall to Analysis 1; however, some
65 differences were large. For example, `weka/NaiveBayes` performed best overall in Analysis 2 but was
66 ranked 28th in Analysis 1.

67



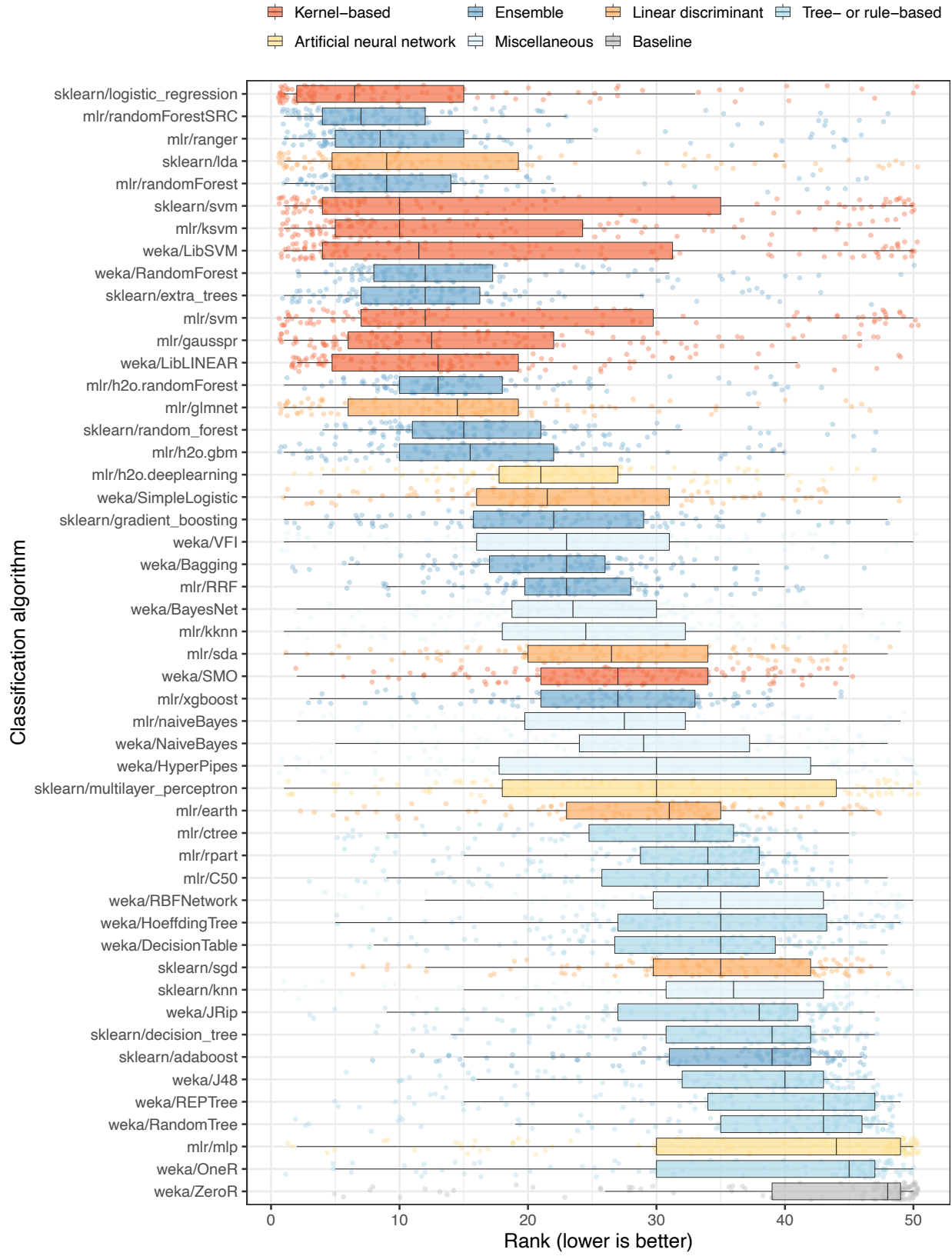
68

69 **Figure S8: Dataset performance by class category when using clinical predictors.** For each class
70 variable across all datasets, we assigned a category representing the type of patient state being predicted.
71 For Analysis 2, we show the predictive performance for each combination of dataset, class variable, and
72 classification algorithm in each class category. We use area under the receiver operating characteristic
73 curve (AUROC) as the metric. The dashed, red line indicates the performance expected by random
74 chance. The top-performing category was “Diagnosis,” which includes class variables associated with a
75 particular disease or subtype. The lowest-performing category was “Patient Characteristic,” which
76 includes variables that indicate whether patients had a family history of cancer, had been diagnosed with
77 multiple tumors, patient performance status, etc.

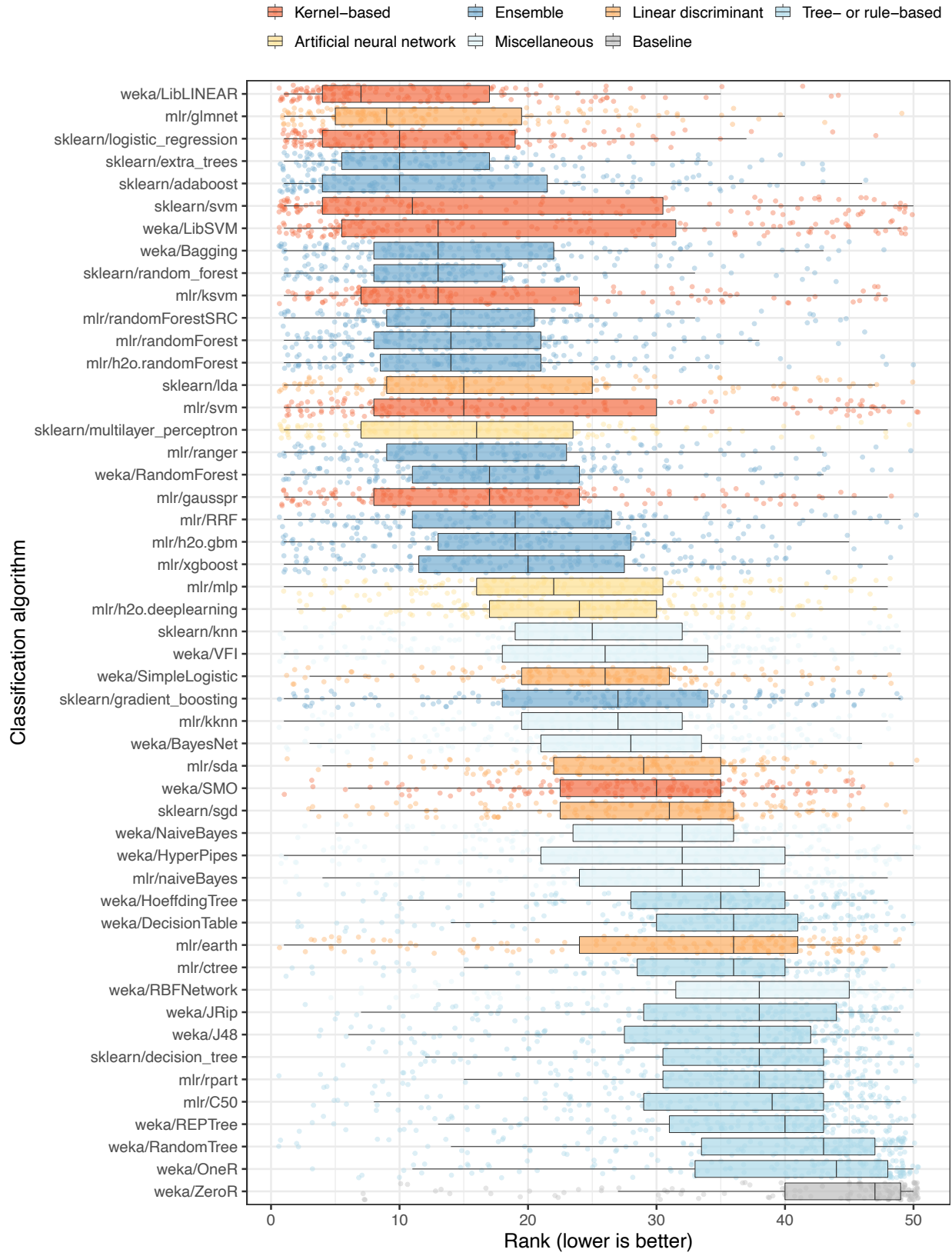


78

79 **Figure S9: Dataset performance by class category when using gene-expression and clinical**
 80 **predictors.** For each class variable across all datasets, we assigned a category representing the type of
 81 patient state being predicted. For Analysis 3, we show the predictive performance for each combination of
 82 dataset, class variable, and classification algorithm in each class category. We use area under the receiver
 83 operating characteristic curve (AUROC) as a metric. The dashed, red line indicates the performance
 84 expected by random chance. As with Analysis 1 (Figure S6), the top-performing category was “Molecular
 85 Marker,” which includes class variables associated with mutation status, immunohistochemistry markers
 86 of protein expression, presence or absence of chromosomal aberrations, etc. The lowest-performing
 87 category was “Patient Characteristic,” which includes variables that indicate whether patients had a
 88 family history of cancer, had been diagnosed with multiple tumors, patient performance status, etc.

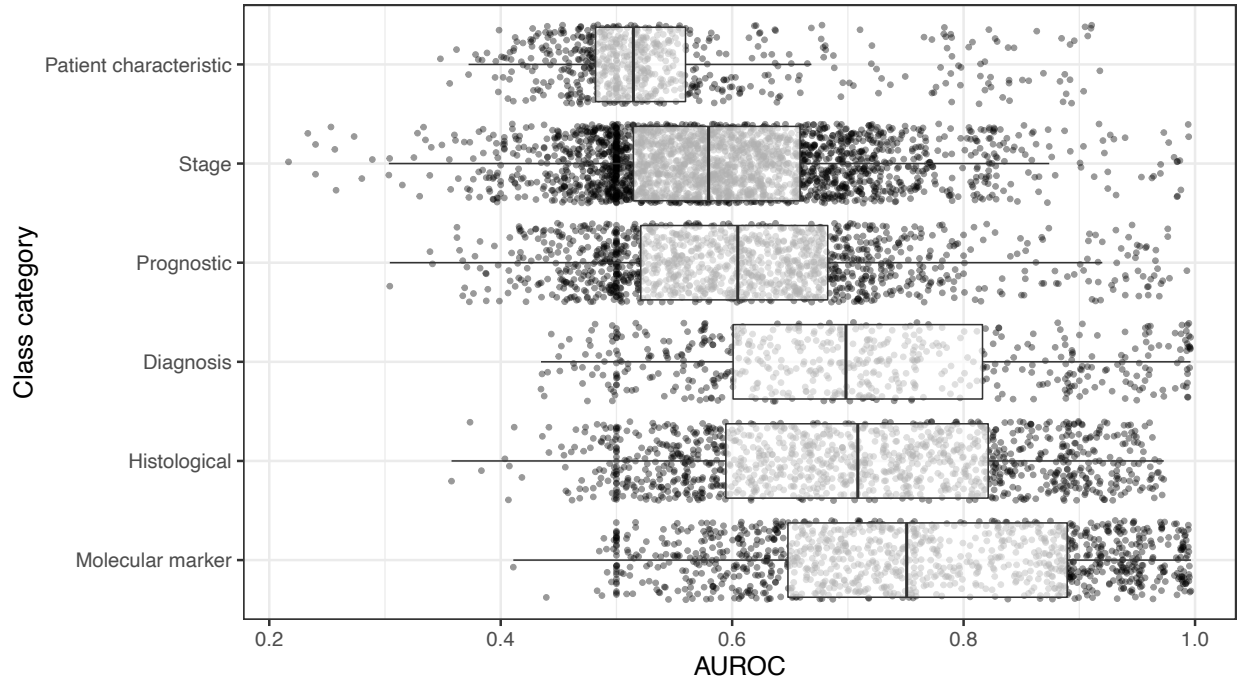


90 **Figure S10: Relative performance of classification algorithms using gene-expression and clinical**
91 **predictors** We predicted patient states using gene-expression and clinical predictors (Analysis 3). For
92 each combination of dataset, class variable, and classification algorithm, we calculated the arithmetic
93 mean of area under the receiver operating characteristic curve (AUROC) values across 50 iterations of
94 Monte Carlo cross-validation. Next we sorted the algorithms based on the average rank across all
95 dataset/class combinations. Each data point that overlays the box plots represents a particular dataset/class
96 combination. The algorithm ranks were highly similar to the ranks for Analysis 1 (Figure S1).



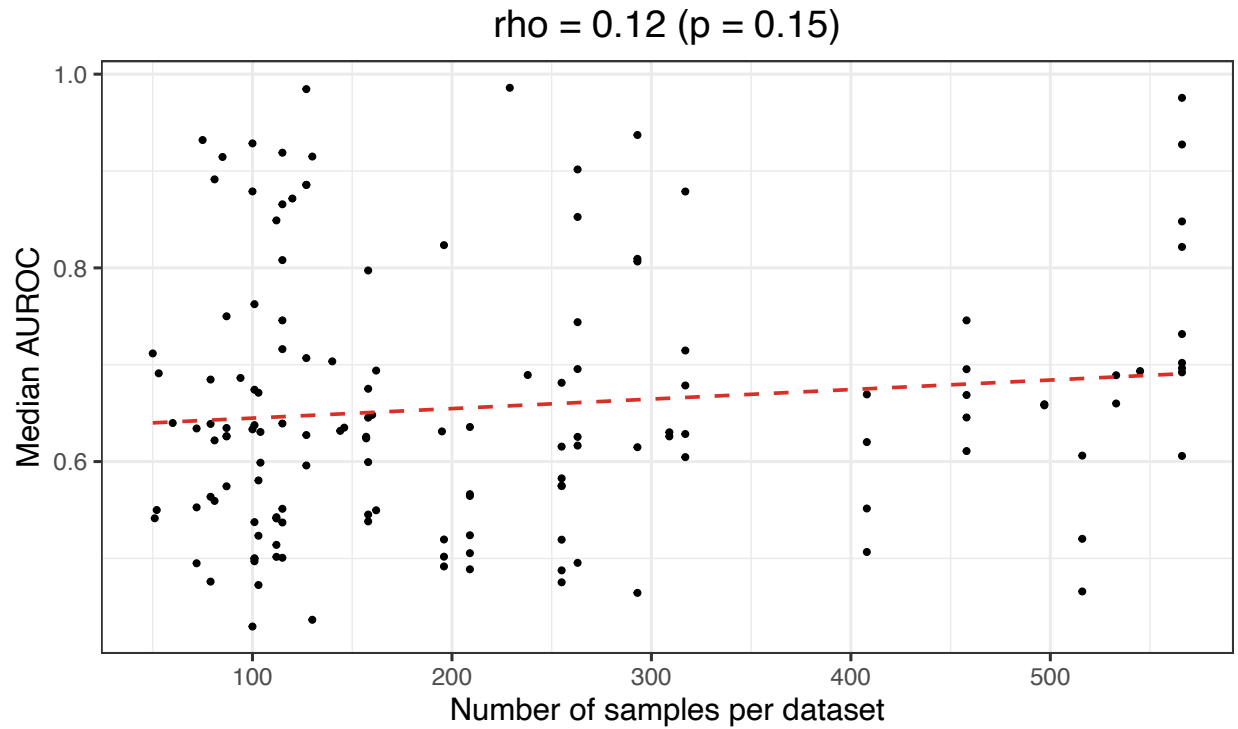
98 **Figure S11: Relative performance of classification algorithms using gene-expression and clinical**
99 **predictors and performing hyperparameter optimization.** We predicted patient states using gene-
100 expression and clinical predictors with hyperparameter optimization (Analysis 4). We used nested cross
101 validation to estimate which hyperparameter combination would be optimal for each algorithm in each
102 training set. For each combination of dataset, class variable, and classification algorithm, we calculated
103 the arithmetic mean of area under the receiver operating characteristic curve (AUROC) values across 5
104 iterations of Monte Carlo cross-validation. Next we sorted the algorithms based on the average rank
105 across all dataset/class combinations. Each data point that overlays the box plots represents a particular
106 dataset/class combination. The algorithm rankings followed similar trends as Analysis 3 (no
107 hyperparameter optimization); however, some differences are notable. For example, the
108 `weka/LibLINEAR` and `mlr/glmnet` algorithms were ranked 11th and 16th in Analysis 3 (Figure
109 S10), but they were ranked 1st and 2nd in this analysis.

110



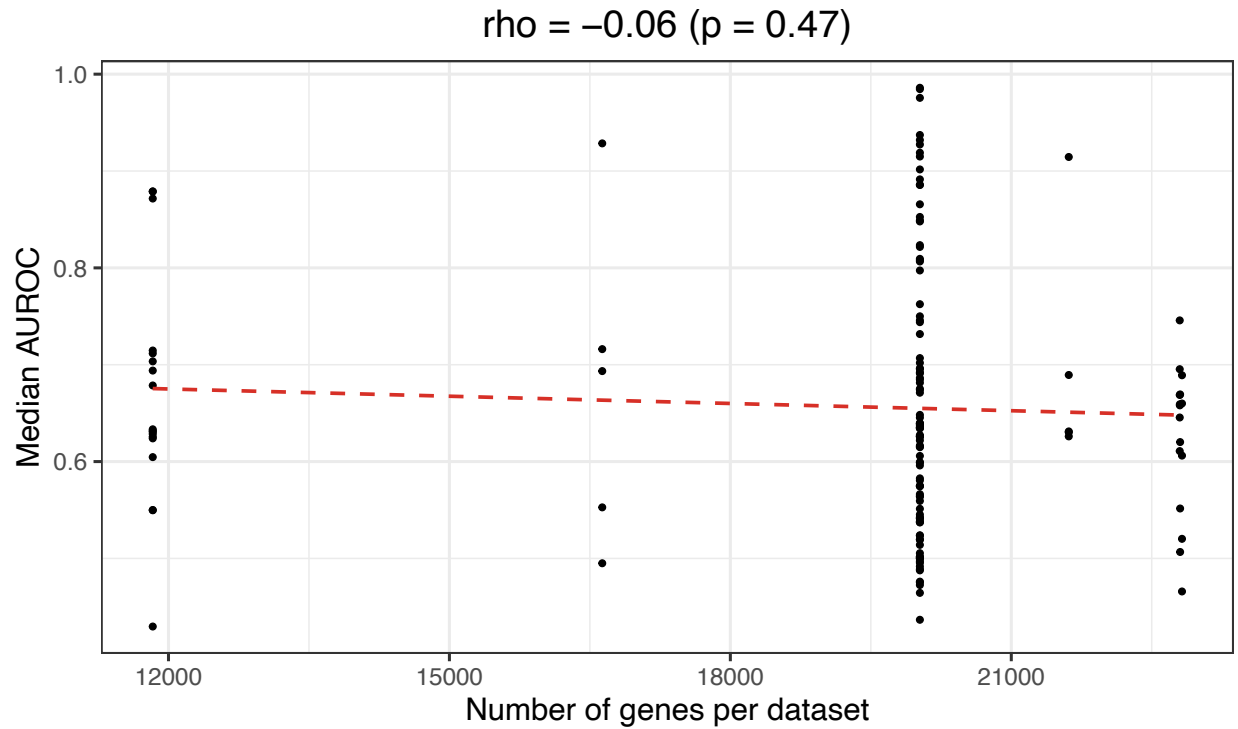
111

112 **Figure S12: Dataset performance by class category when using gene-expression and clinical**
 113 **predictors and performing hyperparameter optimization.** For each class variable across all datasets,
 114 we assigned a category representing the type of patient state being predicted. For Analysis 4, we show the
 115 predictive performance for each combination of dataset, class variable, and classification algorithm in
 116 each class category. We use area under the receiver operating characteristic curve (AUROC) as a metric.
 117 The dashed, red line indicates the performance expected by random chance. The results are similar to
 118 those of Analysis 3 (Figure S9).



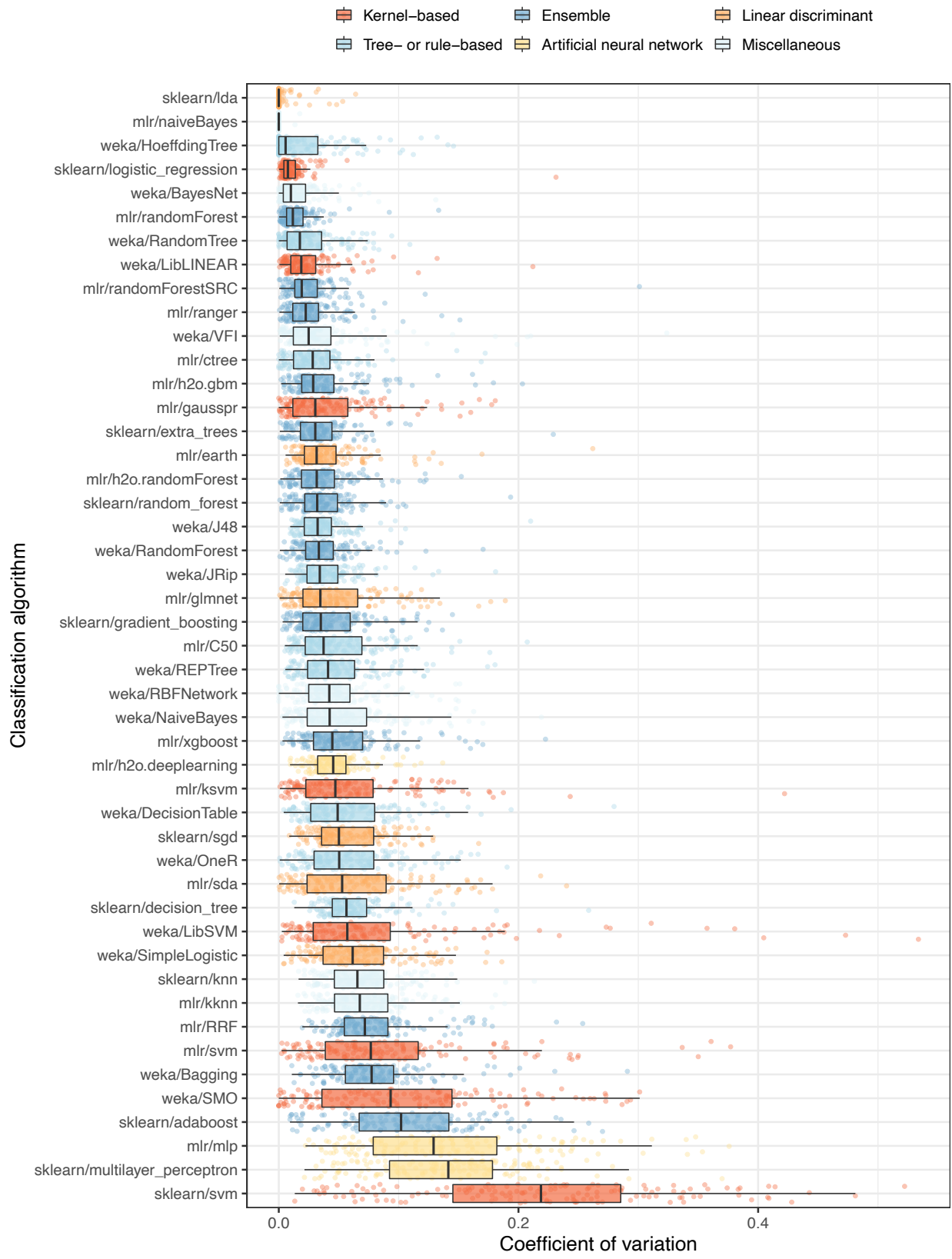
119

120 **Figure S13: Relationship between predictive performance and number of samples per dataset.** The
121 number of patient samples differed by dataset. This scatterplot shows the relationship between the median
122 area under the receiver operating characteristic curve (AUROC) and the number of samples in each
123 dataset. We did not observe a significant relationship between these variables.

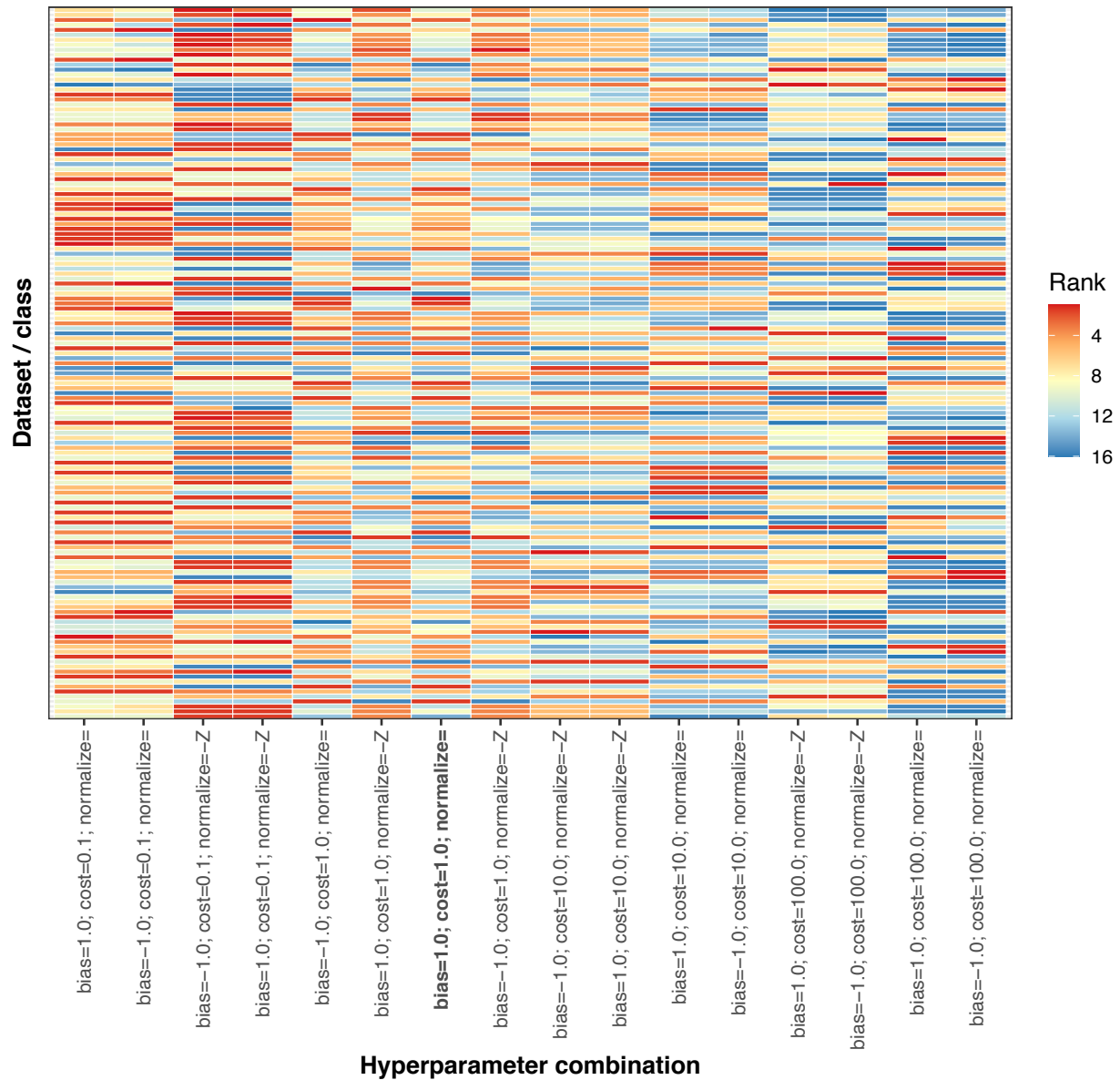


124

125 **Figure S14: Relationship between predictive performance and number of genes per dataset.** Due to
126 differences in gene-expression profiling platforms, we had data for more genes in some datasets than in
127 others. This scatterplot shows the relationship between the median area under the receiver operating
128 characteristic curve (AUROC) and the number of genes in each dataset. We did not observe a significant
129 relationship between these variables.



131 **Figure S15: Variation in predictive performance across hyperparameter combinations.** In Analysis
132 4, we used nested cross validation to evaluate multiple hyperparameter combinations for each
133 classification algorithm. We assessed the extent to which the area under the receiver operating
134 characteristic curve (AUROC) varied across the hyperparameter combinations for each algorithm. For
135 each combination of dataset, class variable, classification algorithm, and hyperparameter set, we averaged
136 AUROC values across 5 Monte Carlo cross-validation iterations. Then we calculated the coefficient of
137 variation for these averaged values across each combination of dataset/class and classification algorithm.
138 Relatively low values indicate that the hyperparameter sets resulted in similar predictive performance. No
139 results are available for 3 algorithms that used only a single hyperparameter option.

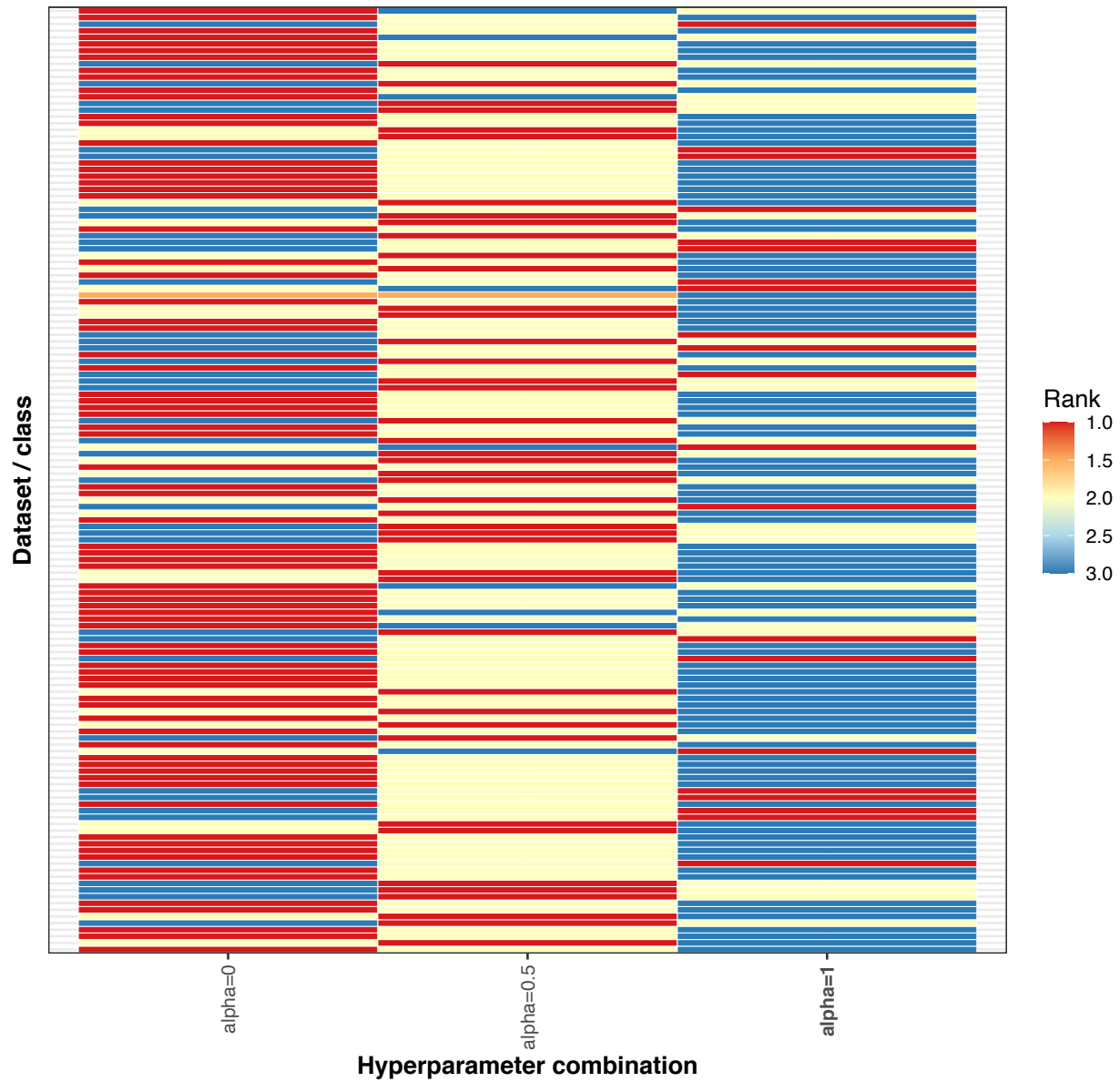


140

141 **Figure S16: Relative performance of different hyperparameter combinations for the**
 142 **weka/LIBLINEAR classification algorithm.** The ShinyLearner software supports 16 hyperparameter
 143 combinations for the weka/LIBLINEAR classification algorithm. In Analysis 4, we used nested cross
 144 validation for hyperparameter optimization. For each combination of dataset and class variable, we
 145 averaged the area under the receiver operating characteristic curve (AUROC) across all (outer) Monte
 146 Carlo cross-validation iterations and then ranked the averages for each hyperparameter combination.

147 Some combinations consistently outperformed other combinations, and the default combination
148 performed suboptimally. Using relatively small `cost` values appeared to improve the performance more
149 than any other option. This hyperparameter controls the regularization strength.

150



151

152 **Figure S17: Relative performance of different hyperparameter combinations for the `mlr/glmnet`**

153 **classification algorithm.** The ShinyLearner software supports 3 hyperparameter combinations for the

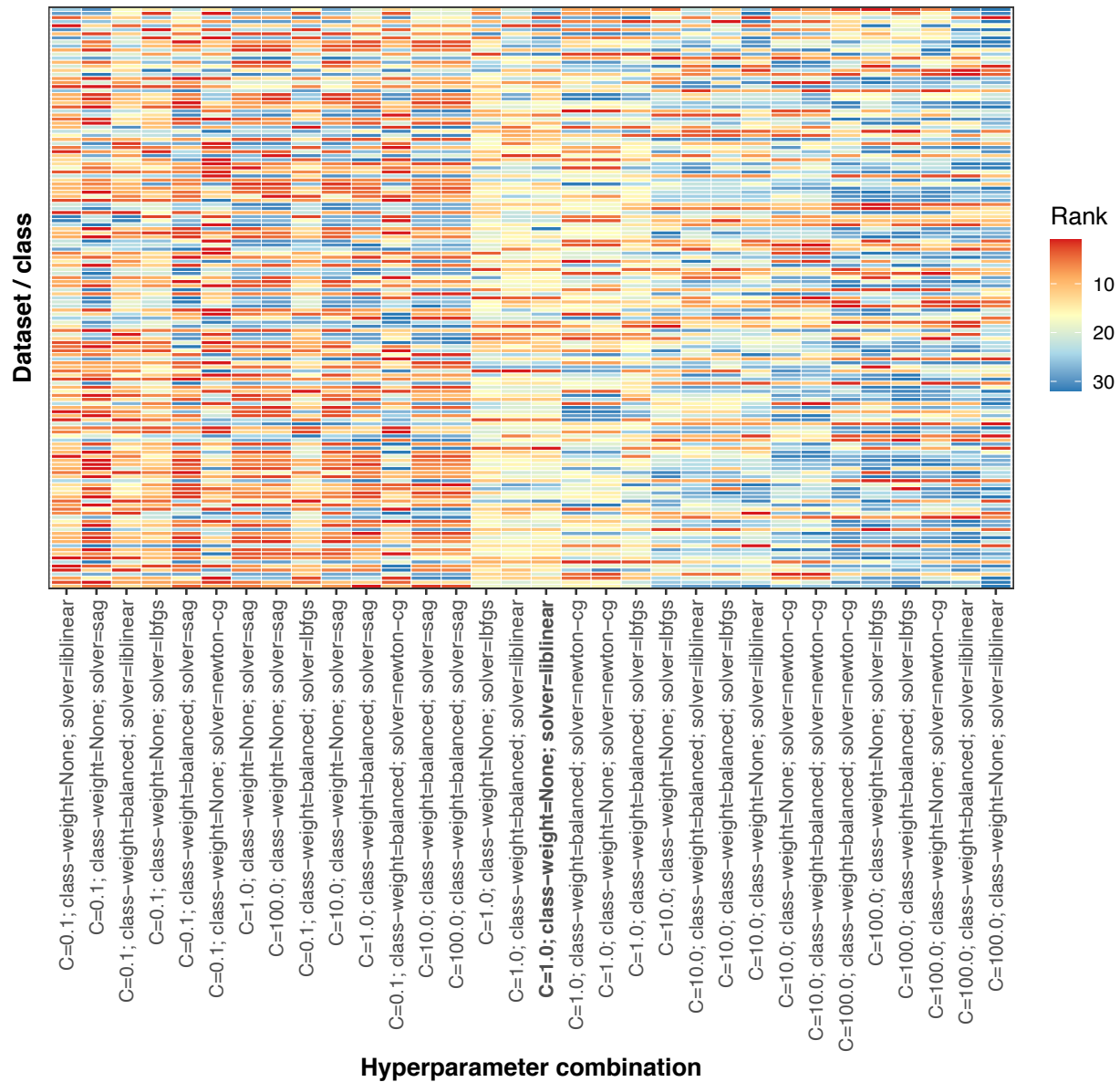
154 `mlr/glmnet` classification algorithm. In Analysis 4, we used nested cross validation for

155 hyperparameter optimization. For each combination of dataset and class variable, we averaged the area

156 under the receiver operating characteristic curve (AUROC) across all (outer) Monte Carlo cross-

157 validation iterations and then ranked the averages for each hyperparameter combination. Using an alpha
158 value of 0.5 or 0 resulted in better performance than a value of 1.

159



160

161

Figure S18: Relative performance of different hyperparameter combinations for the

162

sklearn/logistic_regression classification algorithm. The ShinyLearner software supports

163

32 hyperparameter combinations for the sklearn/logistic_regression classification algorithm.

164

In Analysis 4, we used nested cross validation for hyperparameter optimization. For each combination of

165

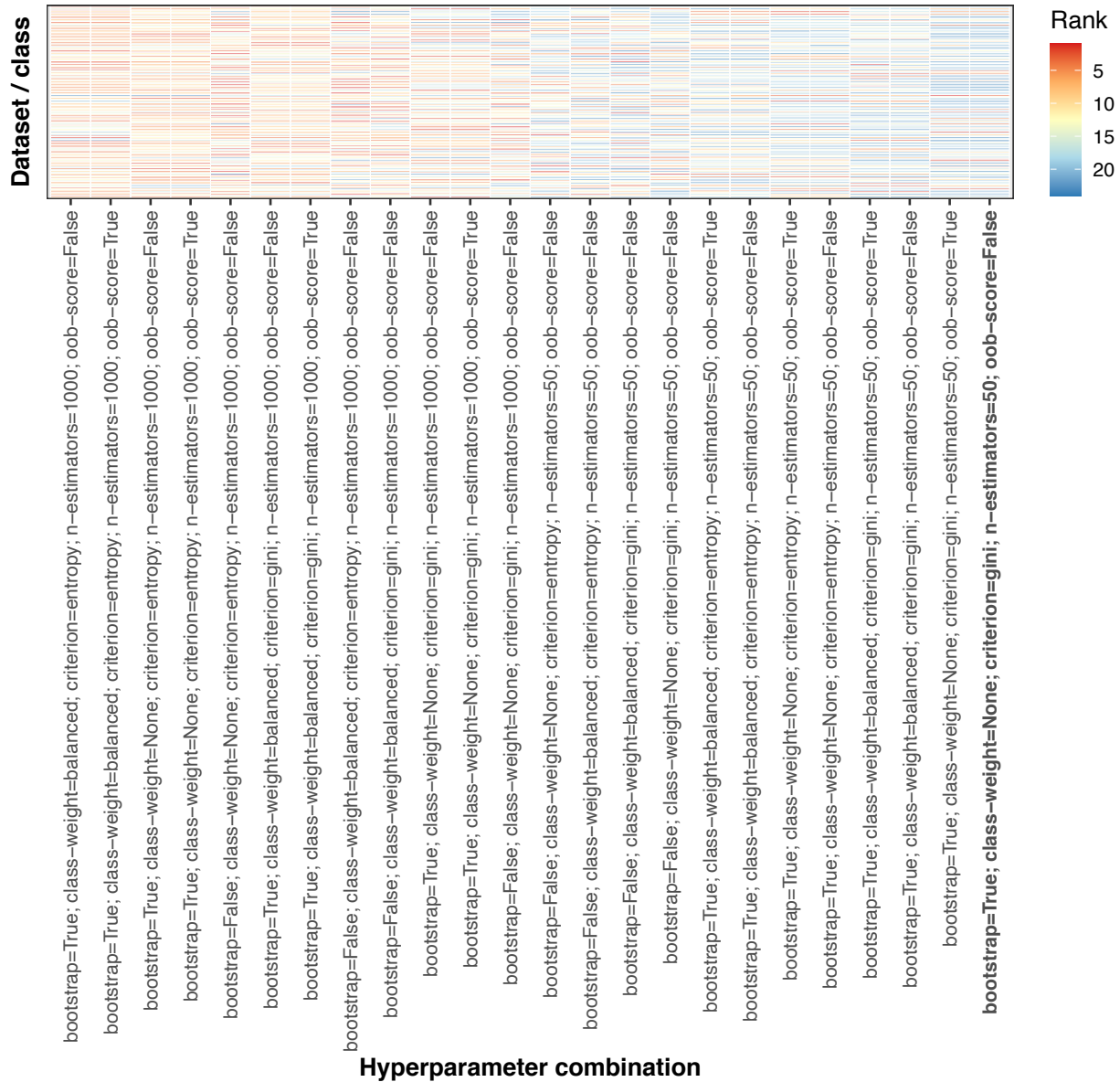
dataset and class variable, we averaged the area under the receiver operating characteristic curve

166

(AUROC) across all (outer) Monte Carlo cross-validation iterations and then ranked the averages for each

167 hyperparameter combination. Some combinations consistently outperformed other combinations, and the
168 default combination performed suboptimally. Using relatively small `cost` values appeared to improve
169 the performance more than any other option. This hyperparameter controls the regularization strength.

170



171

172 **Figure S19: Relative performance of different hyperparameter combinations for the**

173 **sklearn/extra_trees classification algorithm.** The ShinyLearner software supports 24

174 hyperparameter combinations for the sklearn/extra_trees classification algorithm. In Analysis 4,

175 we used nested cross validation for hyperparameter optimization. For each combination of dataset and

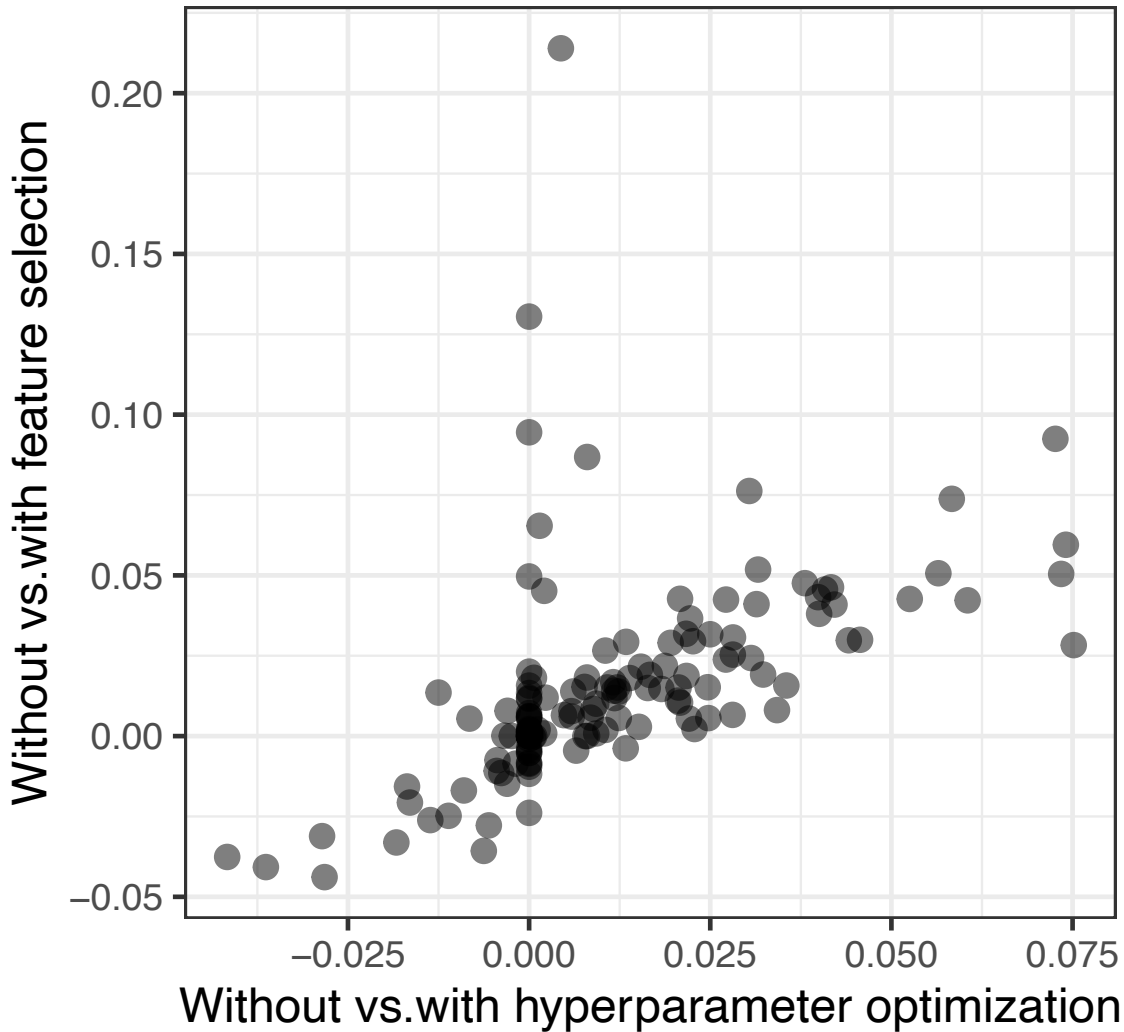
176 class variable, we averaged the area under the receiver operating characteristic curve (AUROC) across all

177 (outer) Monte Carlo cross-validation iterations and then ranked the averages for each hyperparameter

178 combination. Some combinations consistently outperformed other combinations, and the default
179 combination performed suboptimally. Using a larger number ($n = 1000$) of estimators (trees) appeared to
180 improve the performance more than any other option.

181

Spearman's rho = 0.75

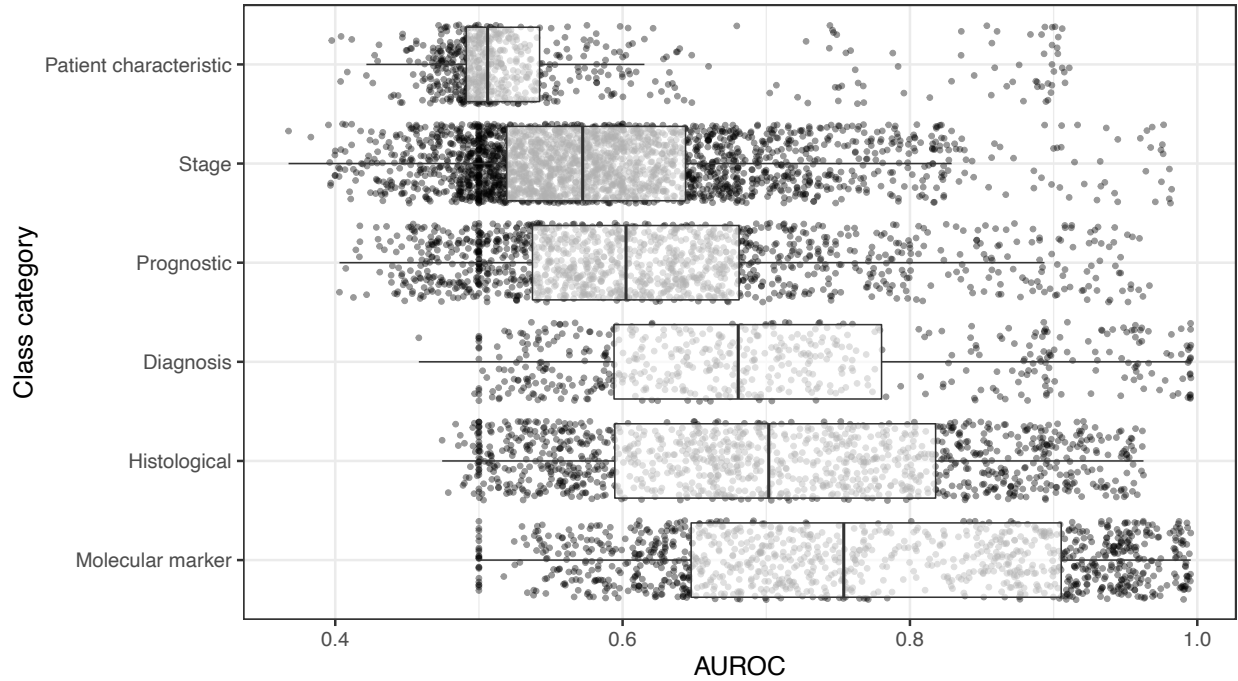


182

183 **Figure S20: Relative predictive performance when using hyperparameter optimization vs. feature**
184 **selection.** We use as a baseline the predictive performance that we attained using default hyperparameters
185 for the classification algorithms (Analysis 3). We quantified predictive performance using the area under
186 the receiver operating characteristic curve (AUROC). This graph shows the increase or decrease in
187 performance when selecting hyperparameters or selecting features relative to the baseline. Each point
188 represents a particular combination of dataset and class variable. Generally, the dataset/class
189 combinations that benefitted from hyperparameter optimization also benefitted from feature selection.

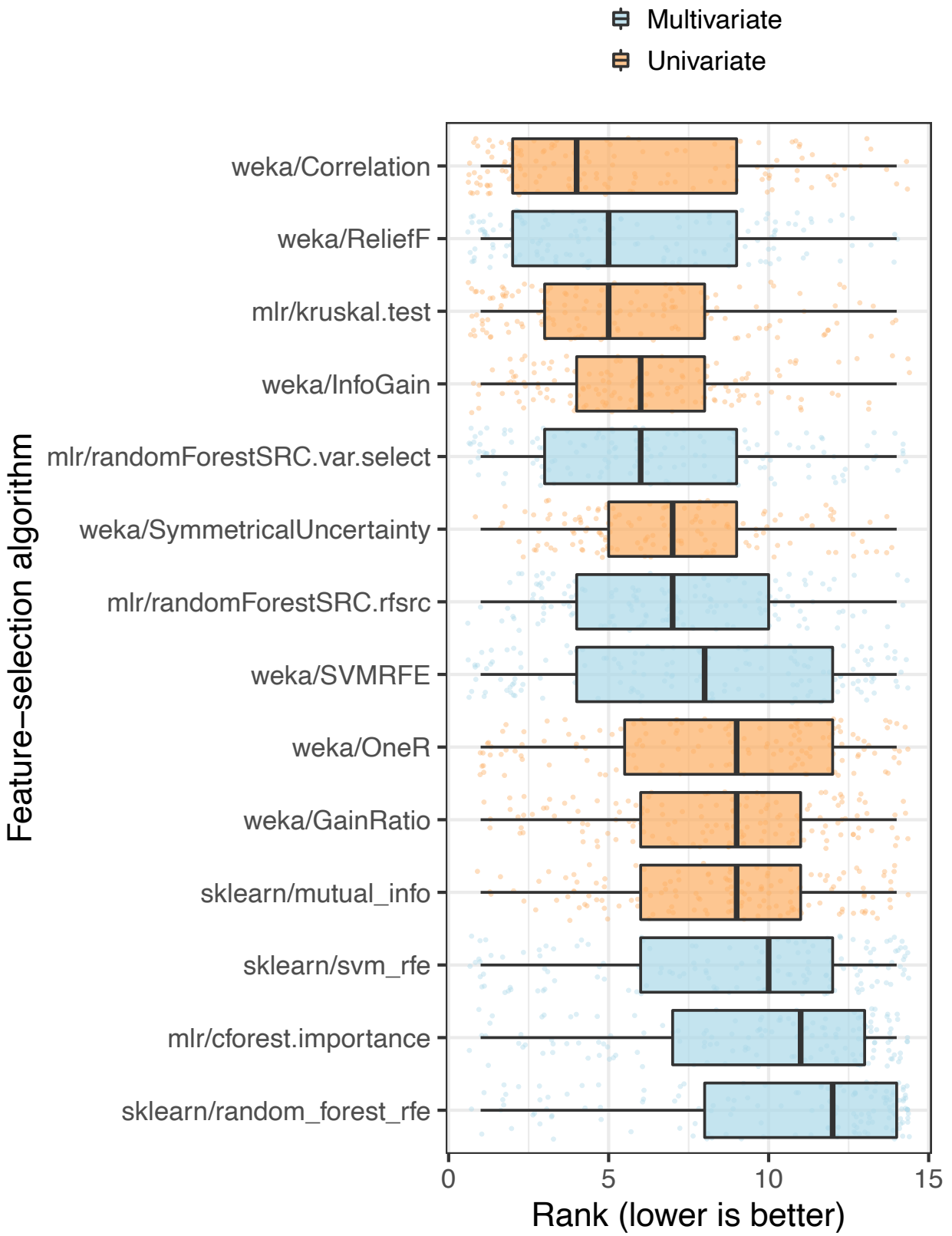
190 However, some dataset/class combinations that did not benefit from hyperparameter optimization did
191 benefit from feature selection.

192

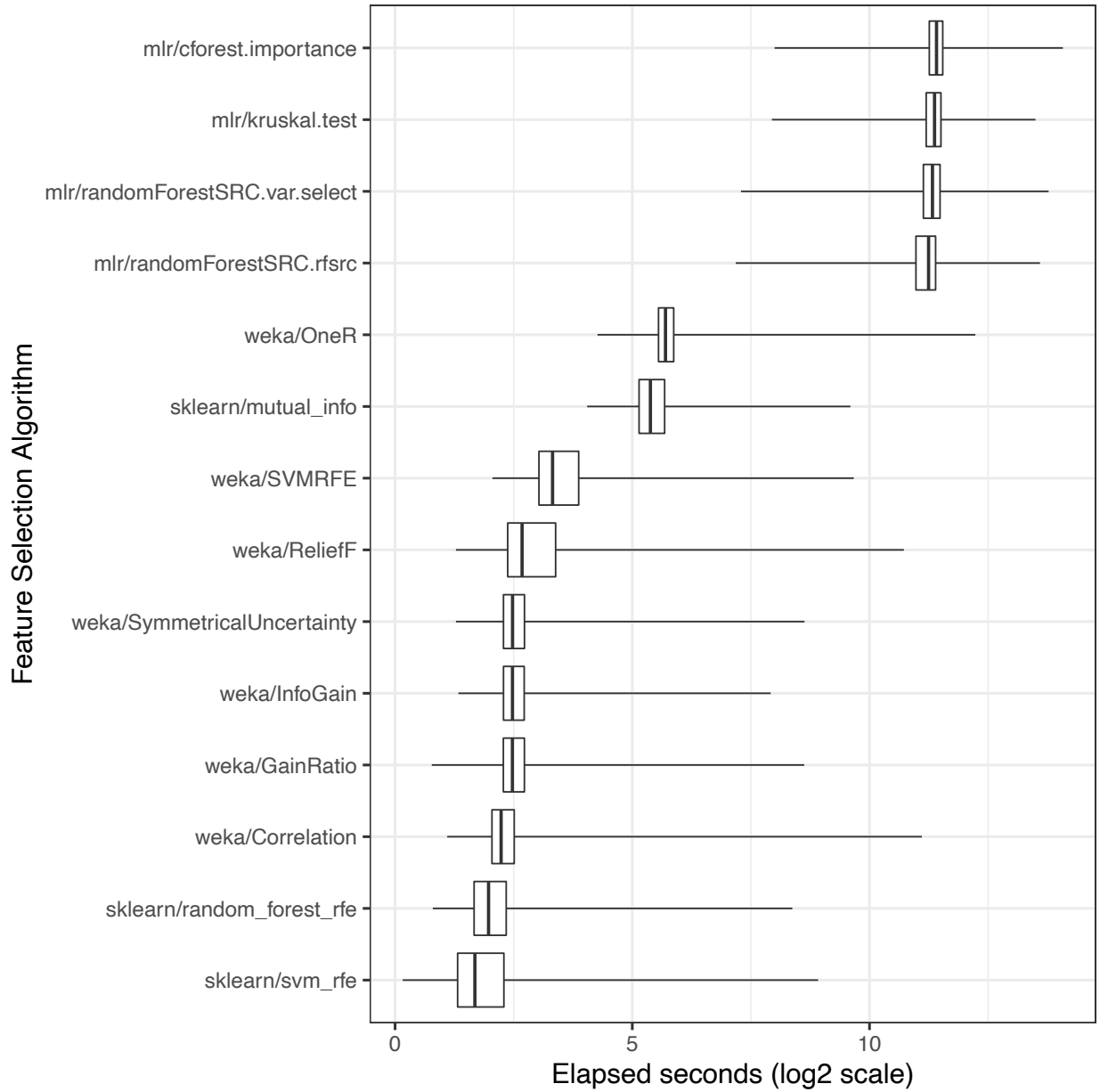


193

194 **Figure S21: Dataset performance by class category when using gene-expression and clinical**
 195 **predictors and performing feature selection.** For each class variable across all datasets, we assigned a
 196 category representing the type of patient state being predicted. For Analysis 5, we show the predictive
 197 performance for each combination of dataset, class variable, and classification algorithm in each class
 198 category. We use area under the receiver operating characteristic curve (AUROC) as a metric. The
 199 dashed, red line indicates the performance expected by random chance. The results are similar to those of
 200 Analyses 3 and 4 (Figures S9, S12).

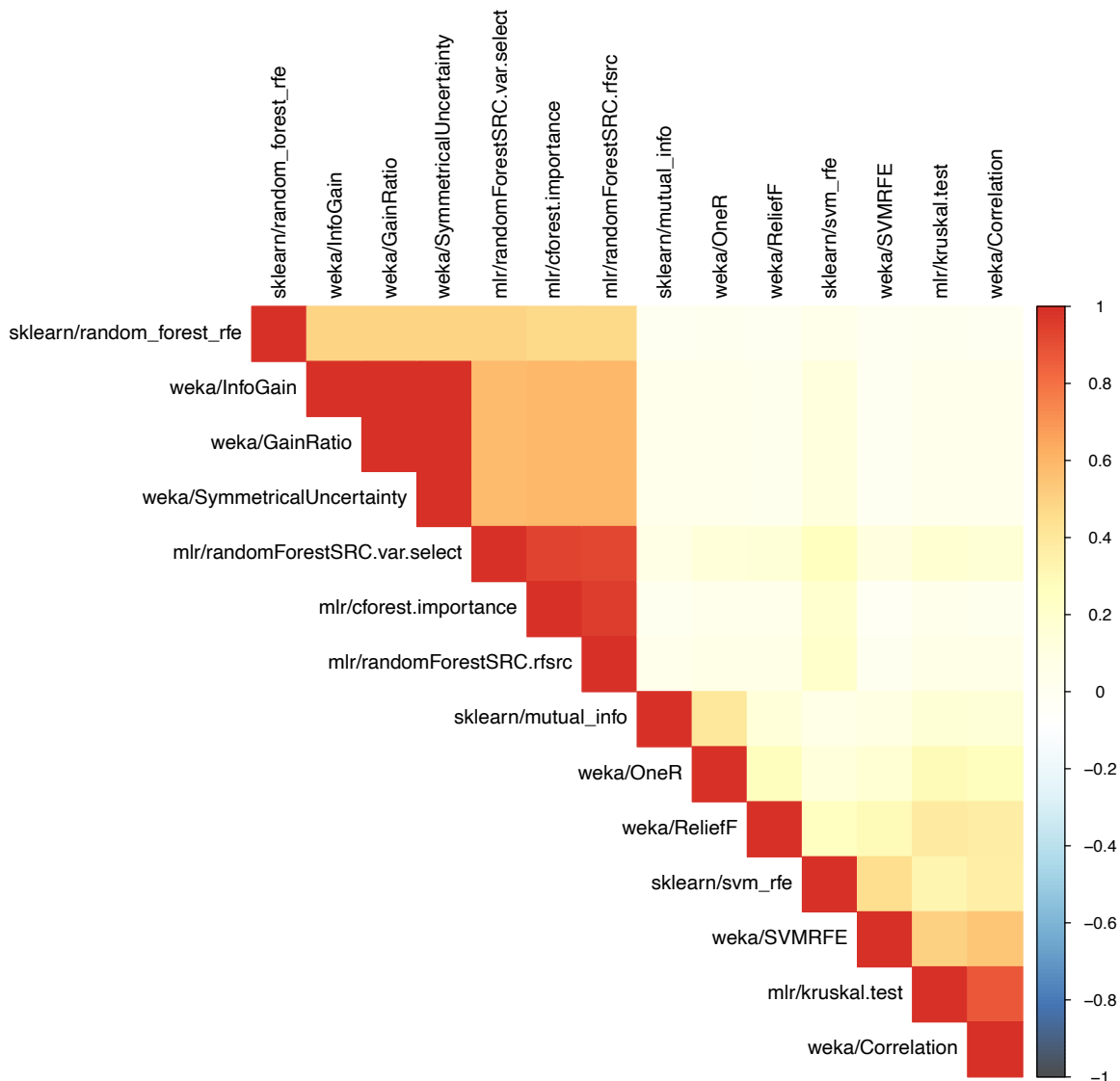


202 **Figure S22: Relative performance of feature-selection algorithms.** For Analysis 5, we used nested
203 cross validation to estimate which features would be most informative for each algorithm in each training
204 set. For each combination of dataset, class variable, and classification algorithm, we ranked the
205 performance of the feature-selection algorithms based on area under the receiver operating characteristic
206 curve (AUROC) and averaged the rankings across 5 iterations of Monte Carlo cross-validation. Each data
207 point that overlays the box plots represents a particular dataset/class combination. Relatively low average
208 ranks are considered optimal. The `weka/Correlation` feature-selection algorithm performed best
209 overall.



210

211 **Figure S23: Execution time per feature-selection algorithm.** In Analysis 5, we used nested cross
 212 validation to estimate which features were most informative for each training set. We calculated the time
 213 (in seconds) required by each feature-selection algorithm to rank the features. Then we averaged these
 214 times across all combinations of dataset, class variable, classification algorithm, and (outer) Monte Carlo
 215 cross-validation iteration. Some feature-selection algorithms were much more computationally intensive
 216 than others.



217

218 **Figure S24: Pairwise correlations of feature ranks between feature-selection algorithms for dataset**

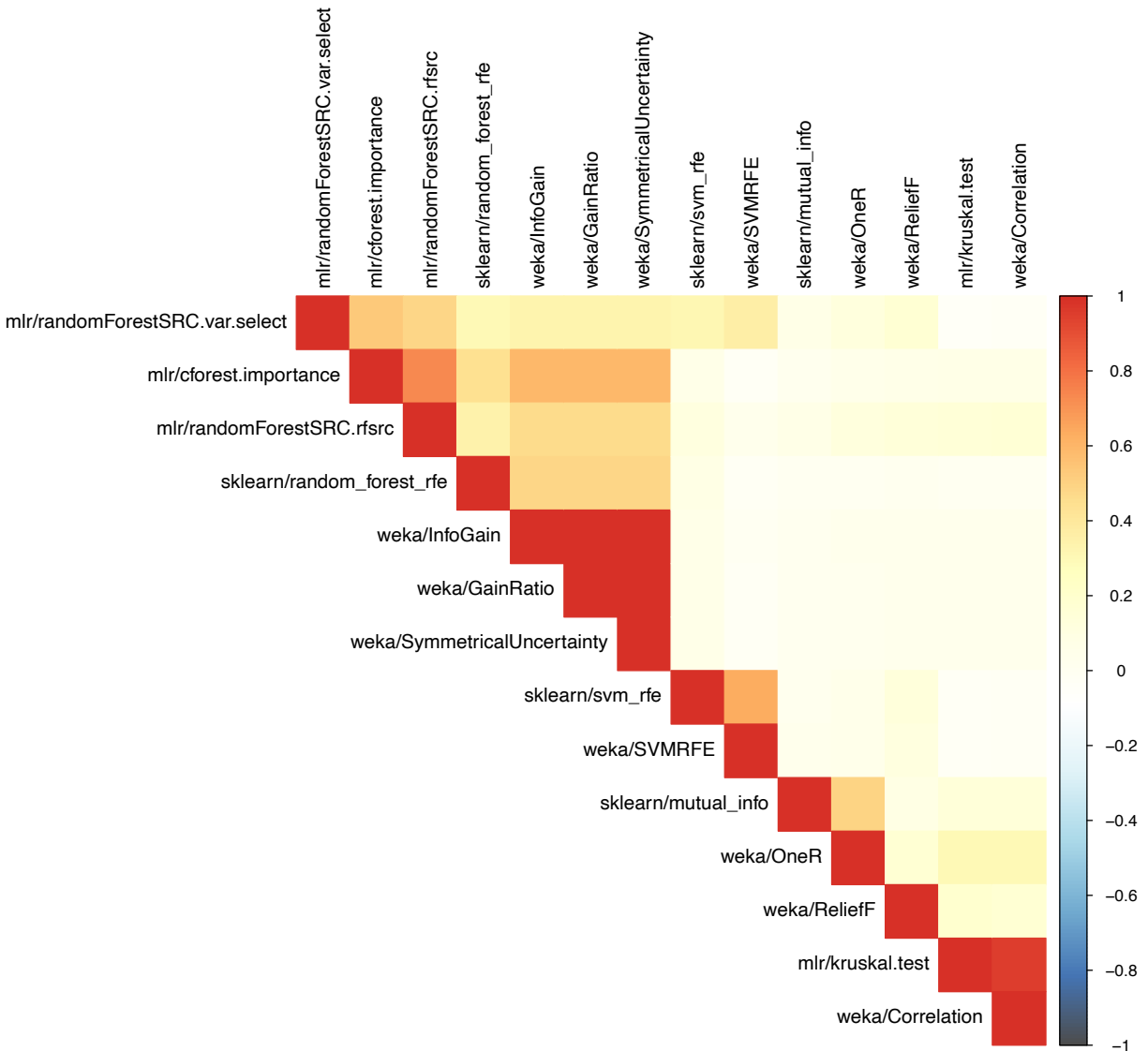
219 **GSE10320.** We used each feature-selection algorithm to rank the genes based on their informativeness

220 for discriminating between relapse and non-relapse outcomes in Wilms tumor patients (GSE10320). After

221 averaging the ranks across cross-validation iterations, we calculated the Spearman correlation coefficient

222 for the feature ranks produced by each pair of algorithms. These coefficients are illustrated as a

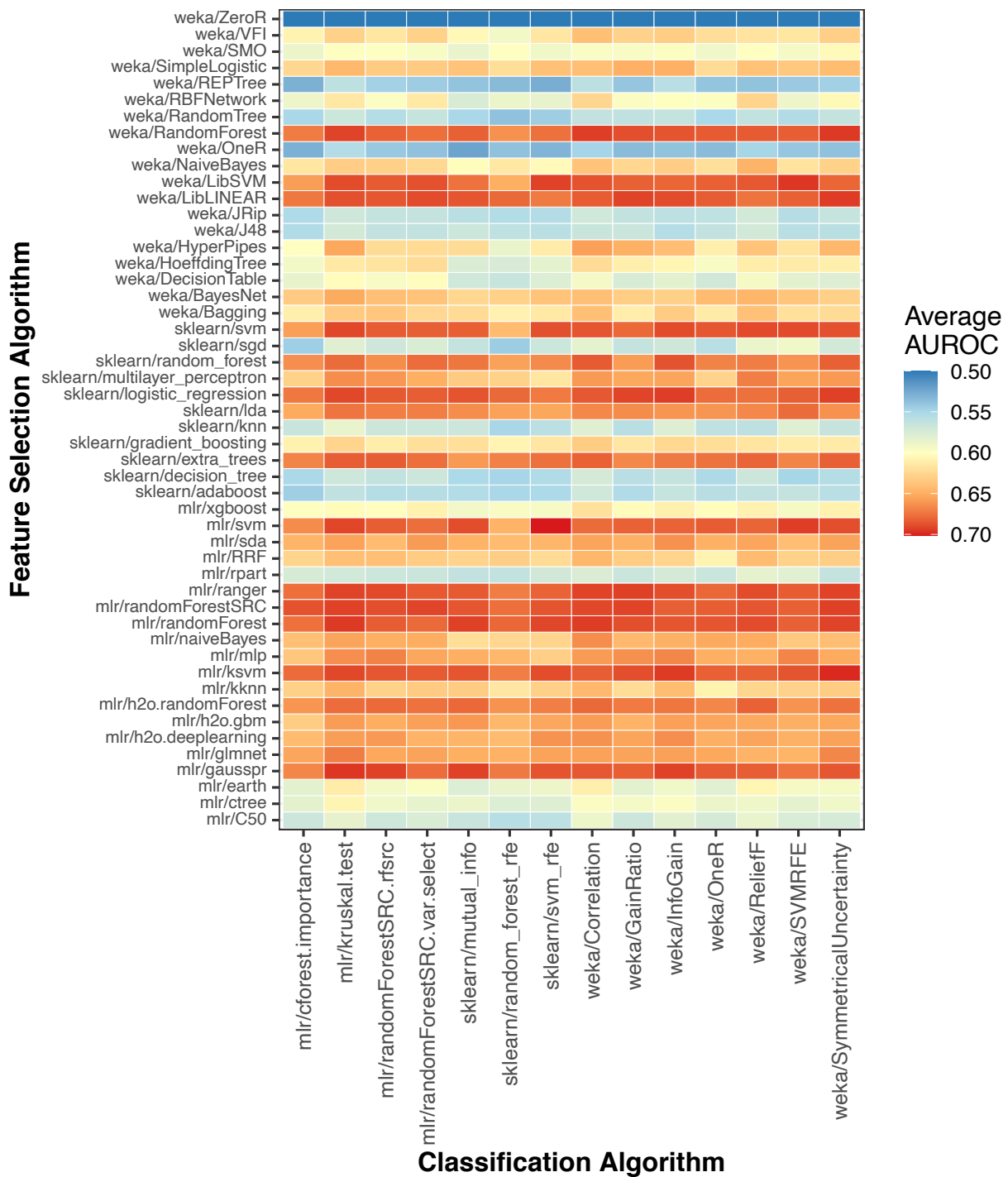
223 correlation plot.



224

225 **Figure S25: Pairwise correlations of feature ranks between feature-selection algorithms for dataset**

226 **GSE46691.** We used each feature-selection algorithm to rank the genes based on their informativeness
 227 for predicting early metastasis following radical prostatectomy (GSE46691). After averaging the ranks
 228 across cross-validation iterations, we calculated the Spearman correlation coefficient for the feature ranks
 229 produced by each pair of algorithms. These coefficients are illustrated as a correlation plot.

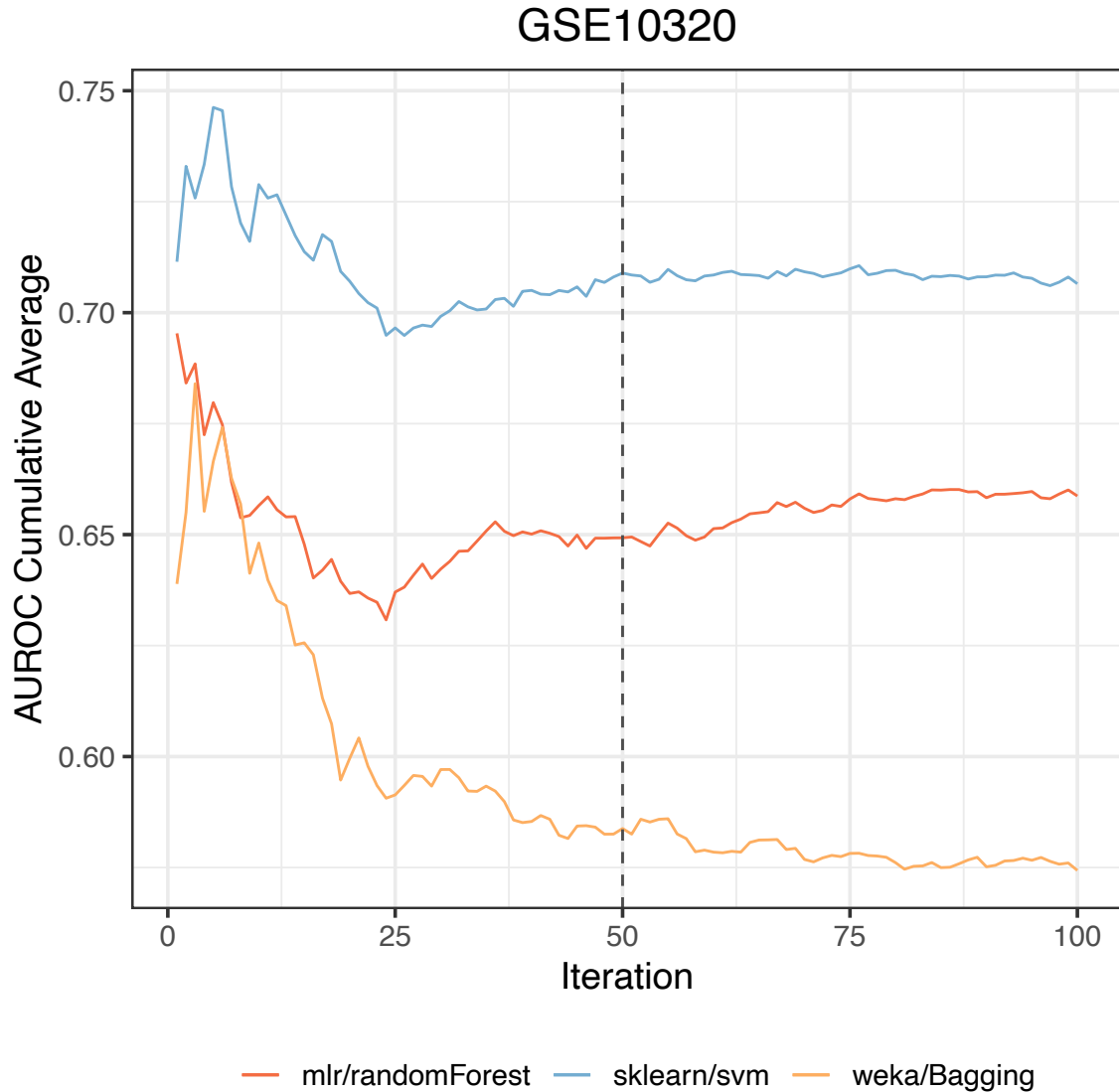


230

231 **Figure S26: Absolute classification performance per combination of feature-selection and**
 232 **classification algorithm.** For each combination of dataset and class variable, we averaged the area under
 233 the receiver operating characteristic curve (AUROC) across all Monte Carlo cross-validation iterations.

234 Then for each combination of feature-selection algorithm and classification algorithm, we calculated the
235 median AUROC across all datasets and class variables.

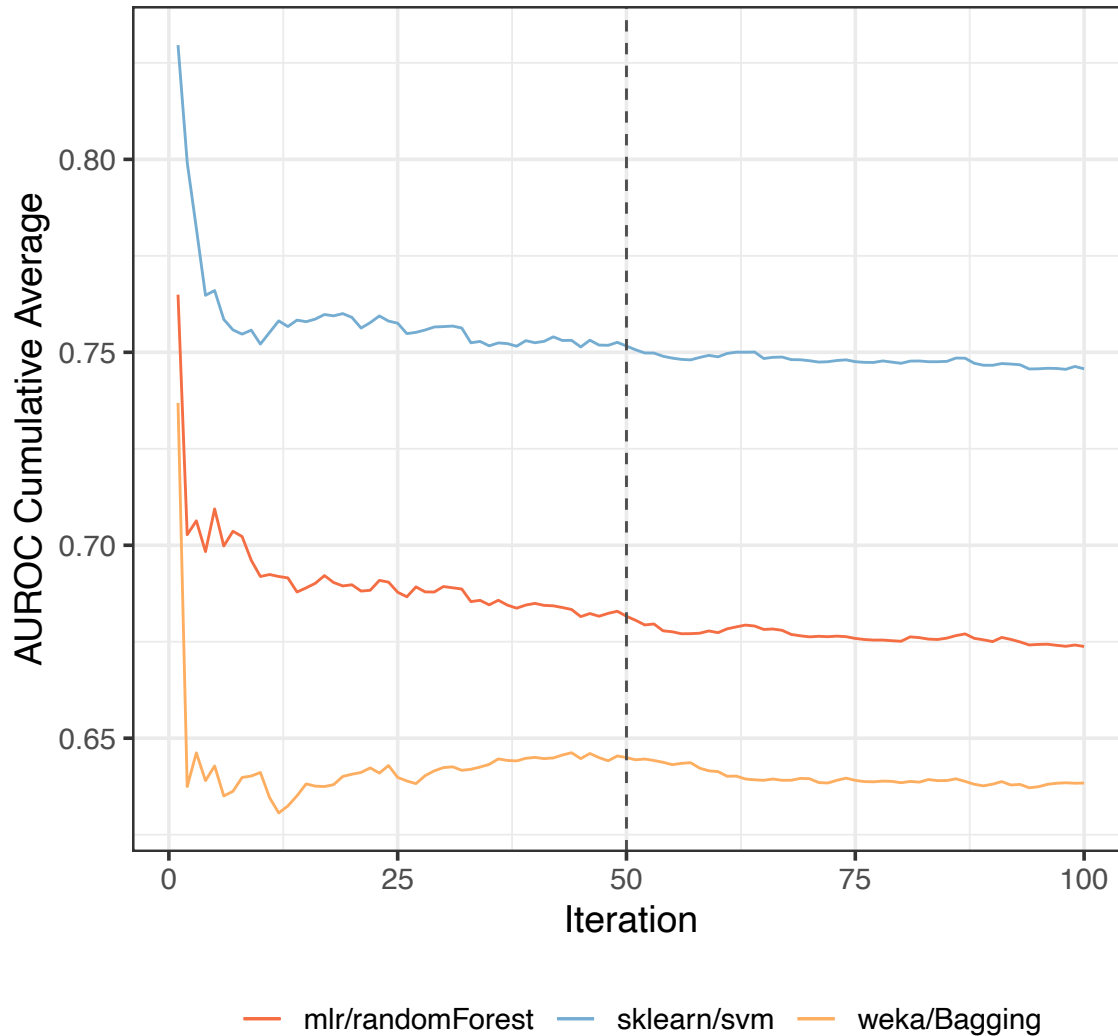
236



237

238 **Figure S27: Stability of classification performance for increasing numbers of cross-validation**
 239 **iterations on dataset GSE10320.** When using gene-expression predictors (Analysis 1), we estimated the
 240 number of Monte Carlo cross-validation iterations that would be sufficient to characterize algorithm
 241 performance. For three classification algorithms, we executed 100 cross-validation iterations on dataset
 242 GSE10320 (predicting relapse vs. non-relapse for Wilms tumor patients). As the number of iterations
 243 increased, we calculated the cumulative average of the area under the receiver operating characteristic
 244 curve (AUROC) for each algorithm. After performing at most 40 iterations, the cumulative averages did
 245 not change more than 0.01 over sequences of 10 iterations.

GSE46691



246

247 **Figure S28: Stability of classification performance for increasing numbers of cross-validation**
248 **iterations on dataset GSE46691.** When using gene-expression predictors (Analysis 1), we estimated the
249 number of Monte Carlo cross-validation iterations that would be sufficient to characterize algorithm
250 performance. For three classification algorithms, we executed 100 cross-validation iterations on dataset
251 GSE46691 (predicting early metastasis following radical prostatectomy). As the number of iterations
252 increased, we calculated the cumulative average of the area under the receiver operating characteristic
253 curve (AUROC) for each algorithm. After performing at most 22 iterations, the cumulative averages did
254 not change more than 0.01 over sequences of 10 iterations.