
FasterRisk: Fast and Accurate Interpretable Risk Scores

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 Over the last century, *risk scores* have been the most popular form of predictive
2 model used in healthcare and criminal justice. Risk scores are sparse linear models
3 with integer coefficients; often these models can be memorized or placed on an
4 index card. Typically, risk scores have been created either without data or by
5 rounding logistic regression coefficients, but these methods do not reliably produce
6 high-quality risk scores. Recent work used mathematical programming, which
7 is computationally slow. We introduce an approach for efficiently producing a
8 collection of high-quality risk scores learned from data. Our approach involves
9 producing a pool of almost-optimal sparse continuous solutions, each with a
10 different support set, using a beam-search algorithm. Each of these continuous
11 solutions is transformed into a separate risk score through a “star search,” where a
12 range of multipliers are considered before rounding the coefficients sequentially
13 to maintain low logistic loss. Our algorithm returns all of these high-quality risk
14 scores for the user to consider. This method completes within minutes and can be
15 impactful in a broad variety of applications.

1 Introduction

17 *Risk scores* are sparse linear models with integer coefficients that predict risks. They are possibly
18 the most popular form of predictive model for high stakes decisions through the last century and are
19 the standard form of model used in criminal justice [4, 20] and medicine [18, 25, 32, 29, 34].
20 Their history dates back to at least the criminal justice work of Burgess [8], where individu-
21 als were assigned integer point scores between 0 and 21 based on their criminal history and
22 demographics that determined their probability of “making good or of failing upon parole.”
23 Other famous risk scores are ar-
24 guably the most widely-used pre-
25 dictive models in healthcare. These
26 include the APGAR score [3], de-
27 veloped in 1952 and given to new-
28 borns, and the CHADS₂ score [17],
29 which estimates stroke risk for
30 atrial fibrillation patients. Figure
31 1 shows an example risk score,
32 which estimates risk of a breast le-
33 sion being malignant.

34 Risk scores have the benefit of be-
35 ing easily memorized; usually their
36 names reveal the full model – for
37 instance, the factors in CHADS₂

1.	Oval Shape	-2 points				
2.	Irregular Shape	4 points	+	...		
3.	Circumscribed Margin	-5 points	+	...		
4.	Spiculated Margin	2 points	+	...		
5.	Age ≥ 60	3 points	+	...		
		SCORE	=			

SCORE	-7	-5	-4	-3	-2	-1
RISK	6.0%	10.6%	13.8%	17.9%	22.8%	28.6%

SCORE	0	1	2	3	4	≥ 5
RISK	35.2%	42.4%	50.0%	57.6%	64.8%	71.4%

Figure 1: Risk score on the mammo dataset [15], whose popula-
tion is biopsy patients. It predicts risk of malignancy of a breast
lesion. Risk score is from FasterRisk on a fold of a 5-CV split.

38 are past Chronic heart failure, Hypertension, Age \geq 75 years, Diabetes, and past Stroke (where past
39 stroke receives 2 points and the others each receive 1 point). For risk scores, counterfactuals are
40 often trivial to compute, even without a calculator. Also, checking that the data and calculations are
41 correct is easier with risk scores than with other approaches. In short, risk scores have been created
42 by humans for a century to support a huge spectrum of applications [2, 21, 28, 40, 41, 43], because
43 humans find them easy to understand.

44 Traditionally, risk scores have been created in two main ways: (1) without data, with expert knowledge
45 only (and validated only afterwards on data), and (2) using a semi-manual process involving manual
46 feature selection and rounding of logistic regression coefficients. That is, these approaches rely
47 heavily on domain expertise and rely little on data. Unfortunately, the converse (relying on data) leads
48 to computationally hard problems: optimizing risk scores over a global objective on data is NP-hard,
49 because in order to produce integer-valued scores, the feasible region must be the integer lattice. There
50 have been only a few approaches to design risk scores automatically [5, 6, 9, 10, 16, 30, 31, 36, 37, 38],
51 but each of these has a flaw that limits its use in practice: the optimization-based approaches use
52 mathematical programming solvers (which require a license) that are slow and scale poorly, and the
53 other methods are randomized greedy algorithms, producing fast but much lower-quality solutions.
54 We need an approach that exhibits the best of both worlds: speed fast enough to operate in a few
55 minutes on a laptop and optimization and search capability as powerful as that of the mathematical
56 programming tools. Our method, FasterRisk, lies at this intersection.

57 One may wonder why simple rounding of ℓ_1 -regularized logistic regression coefficients does not
58 yield sufficiently good risk scores. Past works [35, 37] explain this as follows: the sheer amount of ℓ_1
59 regularization needed to get a very sparse solution leads to large biases and worse loss values, and
60 rounding goes against the performance gradient. For example, consider a set of ℓ_1 coefficients found
61 as [1.45, .87, .83, .47, .23, .15, ...]. This model would be worse than its unregularized counterpart,
62 because of the bias due to the large ℓ_1 term. Its rounded solution is [1,1,1,0,0,0,...], which leads to
63 even worse loss. One could attempt instead to multiply by a constant and then round, but which
64 constant? There are an infinite number of choices. And, even if some value of the multiplier led to
65 minimal loss due to rounding, the bias from the ℓ_1 term still limits the quality of the solution.

66 The algorithm presented here does not have these disadvantages. The steps are: (1) Fast subset search
67 with ℓ_0 optimization (avoiding the bias from ℓ_1). This requires the solution of an NP-hard problem,
68 but our fast subset selection algorithm is able to solve this quickly. We proceed from this accurate
69 sparse continuous solution, preserving both sparseness and accuracy in the next steps. (2) Find a pool
70 of diverse continuous sparse solutions that are almost as good as the solution found in (1) but with
71 different support sets. (3) A “star ray” search, where we search for feasible integer-valued solutions
72 along multipliers of each item in the pool from (2). By using multipliers, the search space resembles
73 a ray of a star because it extends each coefficient in the pool outwards from the origin to search
74 for solutions. To find integer solutions, we perform a local search (a form of sequential rounding).
75 This method yields high performance solutions: we provide a theoretical upper bound on the loss
76 difference between the continuous sparse solution and the rounded integer sparse solution.

77 Through extensive experiments, we show that our proposed method is computationally fast and
78 produces high-quality integer solutions. This work thus provides valuable and novel tools to create
79 risk scores for professionals in many different fields, such as healthcare, finance, and criminal justice.

80 2 Related Work

81 *Optimization-based approaches:* Risk scores, which model $P(y = 1|\mathbf{x})$, are different than threshold
82 classifiers, which predict either $y = 1$ or $y = -1$ given \mathbf{x} . Most work in the area of optimization of
83 integer-valued sparse linear models focuses on classifiers, not risk scores [5, 6, 9, 30, 31, 35, 38, 42].
84 This difference is important, because a classifier generally cannot be calibrated well for use in risk
85 scoring: only its single decision point is optimized. Despite this, several works use the hinge loss
86 to calibrate predictions [6, 9, 30]. All of these optimization-based algorithms use mathematical
87 programming solvers (i.e., integer programming solvers), which tend to be slow and cannot be used
88 on larger problems. However, they can handle both feature selection and integer constraints.

89 To directly optimize risk scores, typically the logistic loss would be used. The RiskSLIM algorithm
90 [37] optimizes the logistic loss regularized with ℓ_0 regularization, subject to integer constraints on the
91 coefficients. RiskSLIM uses callbacks to a MIP solver, alternating between solving linear programs

92 and using branch-and-cut to divide and reduce the search space. The branch-and-cut procedure needs
 93 to keep track of unsolved nodes, whose number increases exponentially with the size of the feature
 94 space. Thus, RiskSLIM’s major challenge is scalability.

95 *Local search-based approaches:* As discussed earlier, a natural way to produce a scoring system or
 96 risk score is by selecting features manually and rounding logistic regression coefficients or hinge-
 97 loss solutions to integers [10, 11, 37]. While rounding is fast, rounding errors discussed earlier
 98 can cause the solution quality to be much worse than that of the optimization-based approaches.
 99 Several works have proposed improvements over traditional rounding. In Randomized Rounding
 100 [10], each coefficient is rounded up or down randomly, based on its continuous coefficient value.
 101 However, randomized rounding does not seem to perform well in practice. Chevaleyre [10] also
 102 proposed Greedy Rounding, where coefficients are rounded sequentially. While this technique
 103 provides theoretical guarantees for greedy rounding for the hinge loss, we have identified a serious
 104 flaw in this argument, rendering the bounds incorrect (see Appendix B). The RiskSLIM paper [37]
 105 proposed SequentialRounding, which, at each iteration, chooses a coefficient to round up or down,
 106 making the best choice according to the regularized logistic loss. This gives better solutions than
 107 other types of rounding, because the coefficients are considered together through their performance
 108 on the loss function, not independently.

109 A drawback of SequentialRounding is that it considers rounding up or down only to the nearest
 110 integer from the continuous solution. By considering *multipliers*, we consider a much larger space
 111 of possible solutions. The idea of multipliers (i.e., “scale and round”) is used for medical scoring
 112 systems [11], though, as far as we know, it has been used only with traditional rounding rather than
 113 SequentialRounding, which could easily lead to poor performance, and we have seen no previous
 114 work that studies how to perform scale-and-round in a systematic, computationally efficient way.
 115 While the general idea of scale-and-round seems simple, it is not: there are an infinite number of
 116 possible multipliers, and, for each one, a number of possible nearby integer coefficient vectors that is
 117 the size of a hypercube, expanding exponentially in the search space.

118 *Sampling Methods:* The Bayesian method of Ertekin et al. [16] samples scoring systems, favoring
 119 those that are simpler and more accurate, according to a prior. “Pooling” [37] creates multiple models
 120 through sampling along the regularization path of ElasticNet. As discussed, when regularization is
 121 tuned high enough to induce sparse solutions, it results in substantial bias and low-quality solutions
 122 (see [35, 37] for numerous experiments on this point). Note that there is a literature on finding diverse
 123 solutions to optimization problems [1], but it only focuses on linear objective functions.

124 **Contributions:** Our contributions include the three-step framework for producing risk scores, the
 125 beam-search based algorithm for logistic regression with bounded coefficients, the search algorithm
 126 to find pools of diverse high-quality continuous solutions, the star search technique using multipliers,
 127 and a theorem guaranteeing the quality of the star search results.

128 **Limitations:** FasterRisk does not provide provably optimal solutions to an NP-hard problem, which
 129 is how it is able to perform in reasonable time for practitioner’s use. FasterRisk’s models should not
 130 be interpreted as causal. FasterRisk creates very sparse generalized additive models and thus has
 131 limited capacity. FasterRisk’s models inherit flaws from data it was trained on. FasterRisk is not yet
 132 customized to a given application, which can be done in future work.

133 3 Methodology

134 Define dataset $\mathcal{D} = \{1, \mathbf{x}_i, y_i\}_{i=1}^n$ (1 is a static feature corresponding to the intercept) and scaled
 135 dataset as $\frac{1}{m} \times \mathcal{D} = \{\frac{1}{m}, \frac{1}{m} \mathbf{x}_i, y_i\}_{i=1}^n$. Our goal is to produce high-quality risk scores within a few
 136 minutes on a small personal computer. We start with an optimization problem similar to RiskSLIM’s
 137 [37], which minimizes the logistic loss subject to sparsity constraints and integer coefficients:

$$\min_{\mathbf{w}, w_0} L(\mathbf{w}, w_0, \mathcal{D}), \quad \text{where } L(\mathbf{w}, w_0, \mathcal{D}) = \sum_{i=1}^n \log(1 + \exp(-y_i(\mathbf{x}_i^T \mathbf{w} + w_0))) \quad (1)$$

such that $\|\mathbf{w}\|_0 \leq k$ and $\mathbf{w} \in \mathbb{Z}^p, \quad \forall j \in [1, \dots, p] \quad w_j \in [-5, 5], \quad w_0 \in \mathbb{Z}.$

138 In practice, the range of these box constraints $[-5, 5]$ is user-defined and can be different for each
 139 coefficient. (We use 5 for ease of exposition.) The $\|\mathbf{w}\|_0$ or integer constraints make the problem
 140 NP-hard, and this is a difficult mixed-integer nonlinear program. Transforming the original features
 141 to all possible dummy variables, as done in other methods [22], changes the model into a (flexible)

Algorithm 1 FasterRisk($\mathcal{D}, k, C, B, \epsilon, T, N_m$) $\rightarrow \{(\mathbf{w}^{+t}, w_0^{+t}, m_t)\}_t$

Input: dataset \mathcal{D} (consisting of feature matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$ and labels $\mathbf{y} \in \mathbb{R}^n$), sparsity constraint k , coefficient constraint $C = 5$, beam search size $B = 10$, tolerance level $\epsilon = 0.3$, number of attempts $T = 50$, number of multipliers to try $N_m = 20$.

Output: a pool P of scoring systems $\{(\mathbf{w}^t, w_0^t), m^t\}$ where t is the index enumerating all found scoring systems with $\|\mathbf{w}^t\|_0 \leq k$ and $\|\mathbf{w}^t\|_\infty \leq C$ and m^t is the corresponding multiplier.

- 1: Call Algorithm [2](#) SparseBeamLR(\mathcal{D}, k, C, B) to find a high-quality solution (\mathbf{w}^*, w_0^*) to the sparse logistic regression problem with continuous coefficients satisfying a box constraint, i.e., solve Problem [1](#). (Algorithm SparseBeamLR will call Algorithm ExpandSuppBy1 as a subroutine, which grows the solution by beam search.)
 - 2: Call Algorithm [5](#) CollectSparseDiversePool($(\mathbf{w}^*, w_0^*), \epsilon, T$), which solves Problem [4](#). Place its output $\{(\mathbf{w}^t, w_0^t)\}_t$ in pool P . $P \leftarrow P \cup \{(\mathbf{w}^t, w_0^t)\}_t$.
 - 3: Send each member t in the pool P , which is (\mathbf{w}^t, w_0^t) , to Algorithm [3](#) StarRaySearch($\mathcal{D}, (\mathbf{w}^t, w_0^t), C, N_m$) to perform a line search among possible multiplier values and obtain an integer solution $(\mathbf{w}^{+t}, w_0^{+t})$ with multiplier m_t . Algorithm [3](#) calls Algorithm [6](#) Auxiliary-LossRounding which conducts the rounding step.
Return the collection of risk scores $\{(\mathbf{w}^{+t}, w_0^{+t}, m_t)\}_t$. If desired, return only the best model according to the logistic loss.
-

142 generalized additive model; even when transformed into risk scores, they can still be as accurate as
143 the best machine learning models [\[37, 39\]](#).

144 To make the solution space substantially larger than $[-5, -4, \dots, 4, 5]^p$, we use *multipliers*. The
145 problem becomes:

$$\min_{\mathbf{w}, w_0, m} L\left(\mathbf{w}, w_0, \frac{1}{m} \mathcal{D}\right), \text{ where } L\left(\mathbf{w}, w_0, \frac{1}{m} \mathcal{D}\right) = \sum_{i=1}^n \log\left(1 + \exp\left(-y_i \frac{\mathbf{x}_i^T \mathbf{w} + w_0}{m}\right)\right) \quad (2)$$

$$\text{such that } \|\mathbf{w}\|_0 \leq k, \mathbf{w} \in \mathbb{Z}^p, \forall j \in [1, \dots, p] w_j \in [-5, 5], w_0 \in \mathbb{Z}, m > 0.$$

146 Note that the use of multipliers does not weaken the interpretability of the risk score: the user still
147 sees integer risk scores comprised of values $w_j \in \{-5, -4, \dots, 4, 5\}$, $w_0 \in \mathbb{Z}$ and points are computed
148 from them. Only the risk conversion table is calculated differently, as $P(Y = 1|\mathbf{x}) = 1/(1 + e^{-f(\mathbf{x})})$
149 where $f(\mathbf{x}) = \frac{1}{m}(\mathbf{w}^T \mathbf{x} + w_0)$.

150 Our method proceeds in three steps, as outlined in Algorithm [1](#). In the first step, it approximately
151 solves the following **sparse logistic regression** problem with a box constraint (but not integer
152 constraints), detailed in Section [3.1](#) and Algorithm [2](#).

$$(\mathbf{w}^*, w_0^*) \in \underset{\mathbf{w}, w_0}{\operatorname{argmin}} L(\mathbf{w}, w_0, \mathcal{D}), \|\mathbf{w}\|_0 \leq k, \mathbf{w} \in \mathbb{R}^p, \forall j \in [1, \dots, p], w_j \in [-5, 5], w_0 \in \mathbb{R}. \quad (3)$$

153 The algorithm gives an accurate and sparse real-valued solution (\mathbf{w}^*, w_0^*) .

154 The second step produces **many near-optimal sparse logistic regression solutions**, again without
155 integer constraints, detailed in Section [3.2](#) and Algorithm [5](#). Algorithm [5](#) uses (\mathbf{w}^*, w_0^*) from the
156 first step to find a set $\{(\mathbf{w}^t, w_0^t)\}_t$ such that for all t and a given threshold ϵ_w :

$$\begin{aligned} (\mathbf{w}^t, w_0^t) \text{ obeys } L(\mathbf{w}^t, w_0^t, \mathcal{D}) &\leq L(\mathbf{w}^*, w_0^*, \mathcal{D}) \times (1 + \epsilon_w) \\ \|\mathbf{w}^t\|_0 &\leq k, \mathbf{w}^t \in \mathbb{R}^p, \forall j \in [1, \dots, p], w_j^t \in [-5, 5], w_0^t \in \mathbb{R}. \end{aligned} \quad (4)$$

157 After these steps, we have a pool of almost-optimal sparse logistic regression models. In the third
158 step, for each coefficient vector in the pool, we **compute a risk score**. It is a feasible integer solution
159 $(\mathbf{w}^{+t}, w_0^{+t})$ to the following, which includes a positive multiplier $m^t > 0$:

$$\begin{aligned} L\left(\mathbf{w}^{+t}, w_0^{+t}, \frac{1}{m^t} \mathcal{D}\right) &\leq L(\mathbf{w}^t, w_0^t, \mathcal{D}) + \epsilon_t, \\ \mathbf{w}^{+t} &\in \mathbb{Z}^p, \forall j \in [1, \dots, p], w_j^{+t} \in [-5, 5], w_0 \in \mathbb{Z}, \end{aligned} \quad (5)$$

160 where we derive a tight theoretical upper bound on ϵ_t . A detailed solution to [\(5\)](#) is shown in Algorithm
161 [6](#) in Appendix [A](#). We do this for a large range of multipliers in Algorithm [3](#). This third step yields a
162 large collection of risk scores, all of which are approximately as accurate as the best sparse logistic
163 regression model that can be obtained. All steps in this process are fast and scalable.

Algorithm 2 SparseBeamLR(\mathcal{D}, k, C, B) $\rightarrow (\mathbf{w}, w_0)$

Input: dataset \mathcal{D} , sparsity constraint k , coefficient constraint C , and beam search size B .

Output: a sparse continuous coefficient vector (\mathbf{w}, w_0) with $\|\mathbf{w}\|_0 = k, \|\mathbf{w}\|_\infty \leq C$.

- 1: Define N_+ and N_- as numbers of positive and negative labels, respectively.
 - 2: $w_0 \leftarrow \log(-N_+/N_-), \mathbf{w} \leftarrow \mathbf{0}$ \triangleright Initialize the intercept and coefficients.
 - 3: $\mathcal{F} \leftarrow \emptyset$ \triangleright Initialize the collection of found supports as an empty set
 - 4: $(\mathcal{W}, \mathcal{F}) \leftarrow \text{ExpandSuppBy1}(\mathcal{D}, (\mathbf{w}, w_0), \mathcal{F}, B)$.
 - 5: **for** $t = 2, \dots, k$ **do** \triangleright Beam search to expand the support
 - 6: $\mathcal{W}_{\text{tmp}} \leftarrow \emptyset$
 - 7: **for** $(\mathbf{w}', w'_0) \in \mathcal{W}$ **do** \triangleright Each of these has support $t - 1$
 - 8: $(\mathcal{W}', \mathcal{F}) \leftarrow \text{ExpandSuppBy1}(\mathcal{D}, (\mathbf{w}', w'_0), \mathcal{F}, B)$. \triangleright Returns $\leq B$ vectors with supp. t .
 - 9: $\mathcal{W}_{\text{tmp}} \leftarrow \mathcal{W}_{\text{tmp}} \cup \mathcal{W}'$
 - 10: **end for**
 - 11: Reset \mathcal{W} to be the B solutions in \mathcal{W}_{tmp} with the smallest logistic loss values.
 - 12: **end for**
 - 13: Pick (\mathbf{w}, w_0) from \mathcal{W} with the smallest logistic loss.
 - 14: **Return** (\mathbf{w}, w_0) .
-

164 3.1 High-quality Sparse Continuous Solution

165 There are many different approaches for sparse logistic regression, including ℓ_1 regularization [33],
166 ElasticNet [44], ℓ_0 regularization [13, 22], orthogonal matching pursuit (OMP) [14, 23], but none
167 of these approaches seem to be able to handle both the box constraints and the sparsity constraint
168 in Problem 3, so we developed a new approach. This approach, in Algorithm 2, SparseBeamLR,
169 uses beam search for best subset selection: each iteration contains several coordinate descent steps
170 to determine whether a new variable should be added to the support, and it clips coefficients to
171 the box $[-5, 5]$ as it proceeds. Hence the algorithm is able to determine, before committing to the
172 new variable, whether it is likely to decrease the loss while obeying the box constraints. This beam
173 search algorithm for solving (3) implicitly uses the assumption that one of the best models of size k
174 implicitly contains variables of one of the best models of size $k - 1$. This type of assumption has
175 been studied in the sparse learning literature [14] (Theorem 5). However, we are not aware of other
176 works applying box constraints or beam search for sparse logistic regression. In Appendix E, we
177 show that our proposed method has higher solution qualities than the OMP method presented in [14].

178 Algorithm 2 calls the ExpandSuppBy1 Algorithm, which has two major steps. The detailed algorithm
179 can be found in Appendix A. For the first step, given a solution \mathbf{w} , we perform optimization on each
180 single coordinate j outside of the current support $\text{supp}(\mathbf{w})$:

$$d_j^* = \underset{d \in [-5, 5]}{\text{argmin}} L(\mathbf{w} + d\mathbf{e}_j, w_0, \mathcal{D}) \text{ for } \forall j \text{ where } w_j = 0. \quad (6)$$

181 We find d_j^* for each j through an iterative thresholding operation, which is done on all coordinates in
182 parallel, iterating several (~ 10) times:

$$\text{for iteration } i: d_j \leftarrow \text{Threshold}(j, d_j, \mathbf{w}, w_0, \mathcal{D}) := \min(\max(c_{d_j}, -5), 5), \quad (7)$$

183 where $c_{d_j} = d_j - \frac{1}{l_j} \nabla_j L(\mathbf{w} + d_j \mathbf{e}_j, w_0, \mathcal{D})$, and l_j is a Lipschitz constant on coordinate j . Importantly,
184 we can perform Equation 7 on all j where $w_j = 0$ in parallel using matrix form.

185 For the second step, after the parallel single coordinate optimization is done, we pick the top B
186 indices (j 's) with the smallest logistic losses $L(\mathbf{w} + d_j^* \mathbf{e}_j)$ and fine tune on the new support:

$$\mathbf{w}_{\text{new}}^j, w_{0\text{new}}^j \in \underset{\mathbf{a} \in [-5, 5]^p, b}{\text{argmin}} L(\mathbf{a}, b, \mathcal{D}) \text{ with } \text{supp}(\mathbf{a}) = \text{supp}(\mathbf{w}) \cup \{j\}. \quad (8)$$

187 This can be done again using a variant of Equation 7 iteratively on all the coordinates in the
188 new support. We get B pairs of $(\mathbf{w}_{\text{new}}^j, w_{0\text{new}}^j)$ through this ExpandSuppBy1 procedure, and the
189 collection of these pairs form the set \mathcal{W}' in Line 8 of Algorithm 2. The ExpandSuppBy1 method is
190 computationally efficient because we are doing parallel single coordinate optimization. This gives
191 the fine-tuning procedure a warm start.

192 **3.2 Collect Sparse Diverse Pool**

193 We now collect the sparse diverse pool. In Section 3.1, our goal was to find a sparse model (\mathbf{w}^*, w_0^*)
 194 with the smallest logistic loss. For high dimensional features or in the presence of highly correlated
 195 features, there could exist many sparse models with almost equally good performance [7]. Let us find
 196 those and turn them into risk scores. We first predefine a tolerance gap level ϵ (usually set to 0.3).
 197 Then, we delete a feature with index j_- in the support $\text{supp}(\mathbf{w}^*)$ and add a new feature with index
 198 j_+ . We select each new index to be j_+ whose logistic loss is within the tolerance gap:

$$\text{Find all } j_+ \text{ s.t. } \min_{a \in [-5, 5]} L(\mathbf{w}^* - w_{j_-}^* \mathbf{e}_{j_-} + a \mathbf{e}_{j_+}, w_0, \mathcal{D}) \leq L(\mathbf{w}^*, w_0^*, \mathcal{D})(1 + \epsilon). \quad (9)$$

199 We fine-tune the coefficients on each of the new supports and then save the new solution in our pool.
 200 Details can be found in Algorithm 5. Swapping one feature at a time is computationally efficient, and
 201 our experiments show it produces sufficiently diverse pools over many datasets.

202 **3.3 “Star” Search for Integer Solutions**

Algorithm 3 StarRaySearch($\mathcal{D}, (\mathbf{w}, w_0), C, N_m$) $\rightarrow (\mathbf{w}^+, w_0^+), m$

Input: dataset \mathcal{D} , a sparse continuous solution (\mathbf{w}, w_0) , coefficient constraint C , and number of
 multipliers to try N_m .

Output: a sparse integer solution (\mathbf{w}^+, w_0^+) with $\|\mathbf{w}^+\|_\infty \leq C$ and multiplier m .

- 1: Define $m_{\max} \leftarrow C / \max|\mathbf{w}|$ as discussed in Section 3.3. If $m_{\max} = 1$, set $m_{\min} \leftarrow 0.5$; if
 $m_{\max} > 1$, set $m_{\min} \leftarrow 1$.
 - 2: Pick N_m equally spaced multiplier values $m_l \in [m_{\min}, m_{\max}]$ for $l \in [1, \dots, N_m]$ and call this
 set $\mathcal{M} = \{m_l\}_l$.
 - 3: Use each multiplier to scale the good continuous solution (\mathbf{w}, w_0) , to obtain $(m_l \mathbf{w}, m_l w_0)$,
 which is a good continuous solution to the rescaled dataset $\frac{1}{m_l} \mathcal{D}$.
 - 4: Send each rescaled solution $(m_l \mathbf{w}, m_l w_0)$ and its rescaled dataset $\frac{1}{m_l} \mathcal{D}$ to Algorithm 6
 AuxiliaryLossRounding($\frac{1}{m_l} \mathcal{D}, m_l \mathbf{w}, m_l w_0$) for rounding. It returns $(\mathbf{w}^{+l}, w_0^{+l}, m_l)$, where
 $(\mathbf{w}^{+l}, w_0^{+l})$ is close to $(m_l \mathbf{w}, m_l w_0)$, and where $(\mathbf{w}^{+l}, w_0^{+l})$ on $\frac{1}{m_l} \mathcal{D}$ has a small logistic loss.
 - 5: Evaluate the logistic loss to pick the best multiplier $l^* \in \text{argmin}_l L(\mathbf{w}^{+l}, w_0^{+l}, \frac{1}{m_l} \mathcal{D})$
 - 6: Return $(\mathbf{w}^{+l^*}, w_0^{+l^*})$ and m_{l^*} .
-

203 The last challenge is how to get an integer solution from a continuous solution. To achieve this, we
 204 use a “star” search that searches along each “ray” of the star, extending each continuous solution
 205 outward from the origin using many values of a multiplier, as shown in Algorithm 3. The star search
 206 provides much more flexibility in finding a good integer solution than simple rounding. The largest
 207 multiplier m_{\max} is set to $5 / \max(|w^*|)$ which will take one of the coefficients to the boundary of the
 208 box constraint at 5. We set the smallest multiplier to be 1.0 and pick N_m (usually 20) equally spaced
 209 points from $[m_{\min}, m_{\max}]$. If $m_{\max} = 1$, we set $m_{\min} = 0.5$ to allow shrinkage of the coefficients.
 210 We scale the coefficients and datasets with each multiplier and round the coefficients to integers using
 211 the sequential rounding technique in Algorithm 6. For each continuous solution (each “ray” of the
 212 “star”), we report the integer solution and multiplier with the smallest logistic loss. This process yields
 213 our collection of risk scores. Note here that a standard line search along the multiplier would not
 214 work because the rounding error is highly non-convex.

215 We briefly discuss how the sequential rounding technique works. Details of this method can be found
 216 in Appendix A. We initialize $\mathbf{w}^+ = \mathbf{w}$. Then we round the fractional part of \mathbf{w}^+ one coordinate at a
 217 time. At each step, some of the w_j^+ ’s are integer-valued (so $w_j^+ - w_j$ is nonzero) and we pick the
 218 coordinate and rounding operation (either floor or ceil) based on which can minimize the following
 219 objective function, where we will round to an integer at coordinate r^* :

$$r^*, v^* \in \underset{r, v}{\text{argmin}} \sum_{i=1}^n l_i^2 \left(x_{ir} (v - w_r) + \sum_{j \neq r} x_{ij} (w_j^+ - w_j) \right)^2, \quad (10)$$

subject to $r \in \{j \mid w_j^+ \notin \mathbb{Z}\}$ and $v \in \{\lfloor w_r^+ \rfloor, \lceil w_r^+ \rceil\}$,

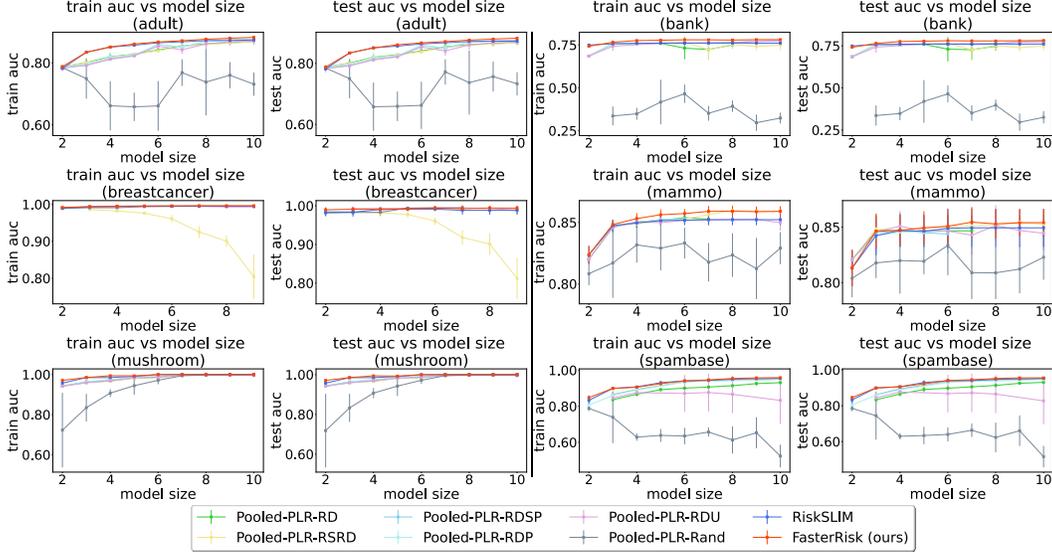


Figure 2: Performance comparison. FasterRisk outperforms baselines due to larger hypothesis space.

220 where l_i is the Lipschitz constant restricted to the rounding interval¹ and can be computed as
 221 $l_i = 1/(1 + \exp(y_i \mathbf{x}_i^T \boldsymbol{\gamma}_i))$ with $\gamma_{ij} = \lfloor w_j \rfloor$ if $y_i x_{ij} > 0$ and $\gamma_{ij} = \lceil w_j \rceil$ otherwise. After we
 222 select r^* and find value v^* , we update \mathbf{w}^+ through $w_{r^*}^+ = v^*$. We repeat this process until \mathbf{w}^+ is
 223 on the integer lattice: $\mathbf{w}^+ \in \mathbb{Z}^p$. The objective function in Equation 10 can be understood as an
 224 auxiliary upper bound of the logistic loss. Our algorithm provides an upper bound on the difference
 225 between the logistic losses of the continuous solution and the final rounded solution before we start
 226 the rounding algorithm (See Theorem C.1). Additionally, during the sequential rounding procedure,
 227 we do not need to perform expensive operations such as logarithms or exponentials as required by the
 228 logistic loss function; the bound and auxiliary function require only sums of squares, not logarithms
 229 or exponentials. Its derivation and proof are in Appendix C

230 **Theorem 3.1.** Let \mathbf{w} be the real-valued coefficients for the logistic regression model with objective
 231 function $L(\mathbf{w}) = \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{x}_i^T \mathbf{w}))$ (the intercept is incorporated). Let \mathbf{w}^+ be the
 232 integer-valued coefficients returned by the AuxiliaryLossRounding method. Furthermore, let $u_j =$
 233 $w_j - \lfloor w_j \rfloor$. Let $l_i = 1/(1 + \exp(y_i \mathbf{x}_i^T \boldsymbol{\gamma}_i))$ with $\gamma_{ij} = \lfloor w_j \rfloor$ if $y_i x_{ij} > 0$ and $\gamma_{ij} = \lceil w_j \rceil$ otherwise.
 234 Then, we have an upper bound on the difference between the loss $L(\mathbf{w})$ and the loss $L(\mathbf{w}^+)$:

$$L(\mathbf{w}^+) - L(\mathbf{w}) \leq \sqrt{n \sum_{i=1}^n \sum_{j=1}^p (l_i x_{ij})^2 u_j (1 - u_j)}. \quad (11)$$

235 **Note.** Our method has a higher prediction capacity than RiskSLIM: its search space is much larger.
 236 Compared to RiskSLIM, our use of the multiplier permits a number of solutions that grows exponen-
 237 tially in k as we increase the multiplier. To see this, consider that for each support of k features, since
 238 logistic loss is convex, it contains a hypersphere in coefficient space. The volume of that hypersphere
 239 is (as usual) $V = \frac{\pi^{k/2}}{\Gamma(\frac{k}{2} + 1)} r^k$ where r is the radius of the hypersphere. If we increase the multiplier to
 240 2, the grid becomes finer by a factor of 2, which is equivalent to increasing the radius by a factor of 2.
 241 Thus, the volume increases by a factor of 2^k . In general, for maximum multiplier m , the search space
 242 is increased by a factor of m^k over RiskSLIM.

243 4 Experiments

244 Our experiments focus on three questions: (1) How good is FasterRisk’s solution quality compared
 245 to baselines? (§4.1) (2) How fast is FasterRisk compared with the state-of-the-art? (§4.2) (3) How

¹The Lipschitz constant here is much smaller than the one in Section 3.1 due to the interval restriction.

246 does each of our proposed technique, including sparse beam search, diverse pool, and multipliers,
 247 contribute to our solution quality? (see Appendix E)

248 We compare with RiskSLIM (the current state-of-the-art), as well as algorithms Pooled-PLR-RD,
 249 Pooled-PLR-RSRD, Pooled-PRL-RDSP, Pooled-PLR-Rand and Pooled-PRL-RDP. These algorithms
 250 were all previously shown to be inferior to RiskSLIM [37]. These methods first find a pool of sparse
 251 continuous solutions using different regularizations of ElasticNet (hence the name ‘‘Pooled Penalized
 252 Logistic Regression’’ – Pooled-PLR) and then round the coefficients with different techniques. Details
 253 are in Appendix D.3. The best solution is chosen from this pool of integer solutions that obeys
 254 the sparsity and box constraints and has the smallest logistic loss. For each dataset, we perform
 255 5-fold cross validation and report training and test AUC. Details about datasets, experimental setup,
 256 evaluation metrics, loss values, and computing platform/environment can be found in Appendix D.
 257 More experimental results appear in Appendix E. Code from [10, 16, 30, 31] is not publicly available.

258 4.1 Solution Quality

259 We first evaluate FasterRisk’s solution quality. Figure 2 shows the training and test AUC on six
 260 datasets (results for training loss appear in Appendix E). FasterRisk (the red line) outperforms all
 261 baselines, consistently obtaining the highest AUC scores on both the training and test sets. Notably,
 262 our method obtains better results than RiskSLIM, which uses a mathematical solver and is the
 263 current state-of-the-art method for scoring systems. This superior performance is due to the use of
 264 multipliers, which increases the complexity of the hypothesis space. A more detailed comparison
 265 between FasterRisk and RiskSLIM appears in Figure 3.

266 FasterRisk performs significantly better than the other baselines for two reasons. First, the continuous
 267 sparse solutions produced by ElasticNet are low quality for very sparse models. Second, it is difficult
 268 to obtain an exact model size by controlling ℓ_1 regularization. For example, Pooled-PLR-RD and
 269 Pooled-PLR-RDSP do not have results for model size 10 on the mammo datasets, because such a
 270 model size does not exist in the pooled solutions after rounding.

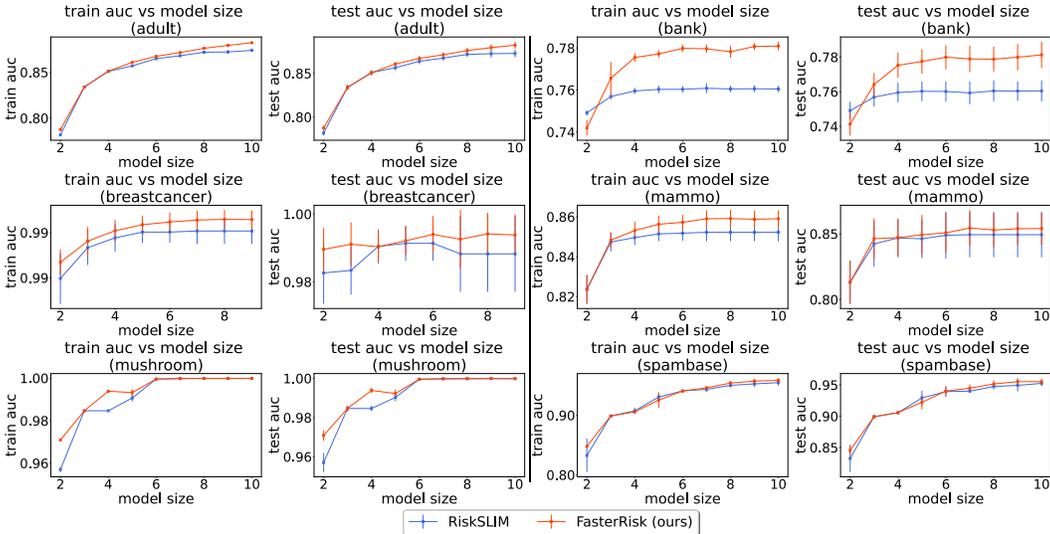


Figure 3: Performance comparison between FasterRisk and RiskSLIM.

271 4.2 Runtime Comparison

272 The major drawback of RiskSLIM is its limited scalability. Figure 4 shows that FasterRisk (red
 273 bars) is significantly faster than RiskSLIM (blue bars) in general. We ran these experiments with a
 274 900 second (15 minute) timeout. RiskSLIM finishes running on small datasets (mammo and breast
 275 cancer), but it times out on the larger datasets, timing out on models larger than 3 features for bank
 276 and spambase, larger than 4 features for adult, and larger than 7 features for mushroom. Thus, we see
 277 that FasterRisk is both faster and more accurate than RiskSLIM.

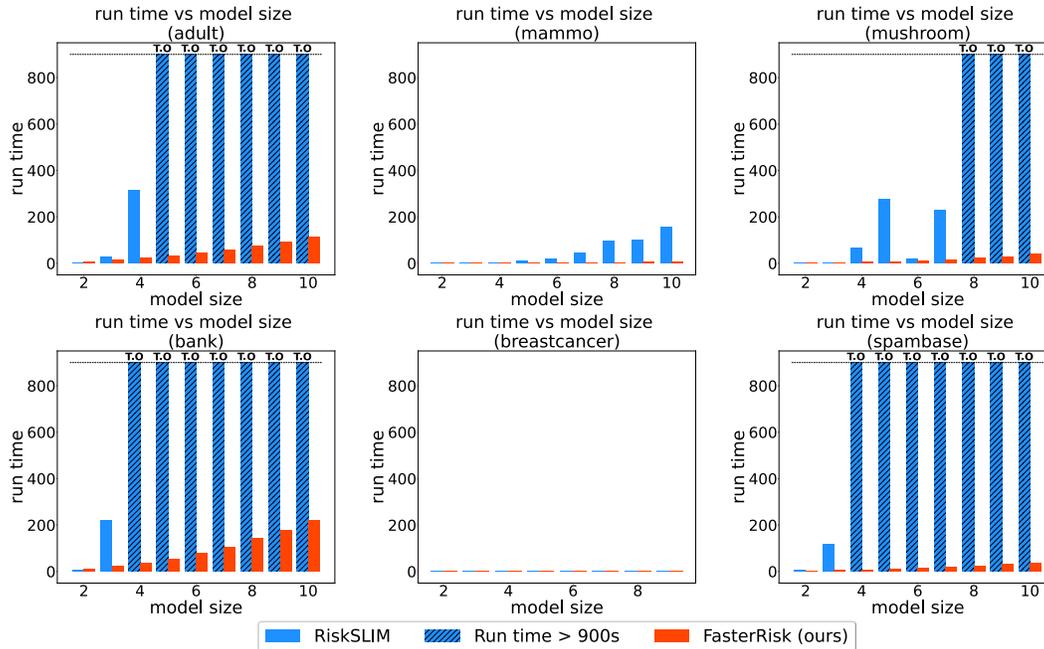


Figure 4: Runtime Comparison. Runtime (in seconds) versus model size for our method FasterRisk (in red) and the RiskSLIM (in blue). The shaded blue bars indicate cases that timed out at 900 seconds. Breastcancer is small and takes approximately 2 seconds for both algorithms.

278 4.3 Example Scoring Systems

279 The main benefit of risk scores is their interpretability. We place a few example risk scores in Table 1 to allow the reader to judge for themselves. More risk scores examples can be found in Appendix F.

1.	no high school diploma	-4 points	...
2.	high school diploma only	-2 points	+ ...
3.	age 22 to 29	-2 points	+ ...
4.	any capital gains	3 points	+ ...
5.	married	4 points	+ ...
SCORE			=

SCORE	<-4	-3	-2	-1	0
RISK	<1.3%	2.4%	4.4%	7.8%	13.6%

SCORE	1	2	3	4	7
RISK	22.5%	35.0%	50.5%	65.0%	92.2%

(a) FasterRisk models for the adult dataset, predicting salary > 50K.

1.	odor=almond	-5 points	...
2.	odor=anise	-5 points	+ ...
3.	odor=none	-5 points	+ ...
4.	odor=foul	5 points	+ ...
5.	gill size=broad	-3 points	+ ...
SCORE			=

SCORE	-8	-5	-3	≥2
RISK	1.62%	26.4%	73.6%	>99.8%

(b) FasterRisk model for the mushroom dataset, predicting whether a mushroom is poisonous.

Table 1: Example FasterRisk models

280

281 5 Conclusion

282 FasterRisk produces a collection of high-quality risk scores within minutes. Its performance owes to
 283 three key ideas: a better algorithm for sparsity and box-constrained continuous models, using a pool
 284 of diverse solutions, and the use of the star search, which leverages multipliers and a new sequential
 285 rounding technique. FasterRisk is suitable for high-stakes decisions, and permits domain experts a
 286 collection of interpretable models to choose from.

- 288 [1] Izuwa Ahanor, Hugh Medal, and Andrew C. Trapp. Diversitree: Computing diverse sets of
289 near-optimal solutions to mixed-integer optimization problems. *arXiv*, 2022.
- 290 [2] Mohamed Farouk Allam. Scoring system for the diagnosis of COVID-19. *The Open Public*
291 *Health Journal*, 13(1), 2020.
- 292 [3] Virginia Apgar. A proposal for a new method of evaluation of the newborn infant. *Current*
293 *Researches in Anesthesia and Analgesia*, 1953(32):260–267, 1953.
- 294 [4] James Austin, Roger Ocker, and Avi Bhati. Kentucky pretrial risk assessment instrument
295 validation. *Bureau of Justice Statistics*, 2010.
- 296 [5] Lieven Billiet, Sabine Van Huffel, and Vanya Van Belle. Interval coded scoring extensions for
297 larger problems. In *Proceedings of the IEEE Symposium on Computers and Communications*,
298 pages 198–203. IEEE, 2017.
- 299 [6] Lieven Billiet, Sabine Van Huffel, and Vanya Van Belle. Interval Coded Scoring: A toolbox for
300 interpretable scoring systems. *PeerJ Computer Science*, 4:e150, 04 2018.
- 301 [7] Leo Breiman. Statistical modeling: The two cultures. *Statistical Science*, 16(3):199–231, 2001.
- 302 [8] Ernest W Burgess. Factors determining success or failure on parole. Illinois Committee on
303 Indeterminate-Sentence Law and Parole Springfield, IL, 1928.
- 304 [9] E. Carrizosa, A. Nogales-Gómez, and D. Romero Morales. Strongly agree or strongly disagree?:
305 Rating features in support vector machines. Technical report, Saïd Business School, University
306 of Oxford, UK, 2013.
- 307 [10] Yann Chevaleyre, Frédéric Kriche, and Jean-Daniel Zucker. Rounding methods for discrete
308 linear classification. In *International Conference on Machine Learning*, pages 651–659. PMLR,
309 2013.
- 310 [11] TJ Cole. Algorithm as 281: scaling and rounding regression coefficients to integers. *Applied*
311 *statistics*, pages 261–268, 1993.
- 312 [12] Lorrie Faith Cranor and Brian A LaMacchia. Spam! *Communications of the ACM*, 41(8):74–83,
313 1998.
- 314 [13] Antoine Dedieu, Hussein Hazimeh, and Rahul Mazumder. Learning sparse classifiers: Continu-
315 ous and mixed integer optimization perspectives. *arXiv preprint arXiv:2001.06471*, 2020.
- 316 [14] Ethan R Elenberg, Rajiv Khanna, Alexandros G Dimakis, and Sahand Negahban. Restricted
317 strong convexity implies weak submodularity. *The Annals of Statistics*, 46(6B):3539–3568,
318 2018.
- 319 [15] Matthias Elter, Rüdiger Schulz-Wendtlund, and Thomas Wittenberg. The prediction of breast
320 cancer biopsy outcomes using two cad approaches that both emphasize an intelligible decision
321 process. *Medical physics*, 34(11):4164–4172, 2007.
- 322 [16] Sÿda Ertekin and Cynthia Rudin. A bayesian approach to learning scoring systems. *Big Data*,
323 2015.
- 324 [17] Brian F Gage, Amy D Waterman, William Shannon, Michael Boechler, Michael W Rich, and
325 Martha J Radford. Validation of clinical classification schemes for predicting stroke. *The*
326 *Journal of the American Medical Association*, 285(22):2864–2870, 2001.
- 327 [18] Ronald C Kessler, Lenard Adler, Minnie Ames, Olga Demler, Steve Faraone, EVA Hiripi,
328 Mary J Howes, Robert Jin, Kristina Secnik, Thomas Spencer, and et al. The world health
329 organization adult ADHD self-report scale (ASRS): a short screening scale for use in the general
330 population. *Psychological Medicine*, 35(02):245–256, 2005.
- 331 [19] Ron Kohavi et al. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In
332 *Kdd*, volume 96, pages 202–207, 1996.
- 333 [20] Edward Latessa, Paula Smith, Richard Lemke, Matthew Makarios, and Christopher Lowenkamp.
334 Creation and validation of the ohio risk assessment system: Final report. *Center for Criminal*
335 *Justice Research, School of Criminal Justice, University of Cincinnati, Cincinnati, OH. Retrieved*
336 *from http://www.ocjs.ohio.gov/ORAS_FinalReport.pdf*, 2009.
- 337 [21] Ji Yeon Lee, Byung-Ho Nam, Mhinjine Kim, Jongmin Hwang, Jin Young Kim, Miri Hyun,
338 Hyun Ah Kim, and Chi-Heum Cho. A risk scoring system to predict progression to severe
339 pneumonia in patients with COVID-19. *Scientific reports*, 12(1):1–8, 2022.
- 340 [22] Jiachang Liu, Chudi Zhong, Margo Seltzer, and Cynthia Rudin. Fast sparse classification for
341 generalized linear and additive models. In *Proceedings of Artificial Intelligence and Statistics*
342 *(AISTATS)*, 2022.
- 343 [23] Aurelie Lozano, Grzegorz Swirszcz, and Naoki Abe. Group orthogonal matching pursuit for
344 logistic regression. In *Proceedings of the fourteenth international conference on artificial*

- 345 *intelligence and statistics*, pages 452–460. JMLR Workshop and Conference Proceedings, 2011.
- 346 [24] Olvi L Mangasarian, W Nick Street, and William H Wolberg. Breast cancer diagnosis and
347 prognosis via linear programming. *Operations Research*, 43(4):570–577, 1995.
- 348 [25] Rui P Moreno, Philipp GH Metnitz, Eduardo Almeida, Barbara Jordan, Peter Bauer, Ri-
349 cardo Abizanda Campos, Gaetano Iapichino, David Edbrooke, Maurizia Capuzzo, and Jean-
350 Roger Le Gall. SAPS 3 - from evaluation of the patient to evaluation of the intensive care unit.
351 part 2: Development of a prognostic model for hospital mortality at ICU admission. *Intensive
352 Care Medicine*, 31(10):1345–1355, 2005.
- 353 [26] Sérgio Moro, Paulo Cortez, and Paulo Rita. A data-driven approach to predict the success of
354 bank telemarketing. *Decision Support Systems*, 62:22–31, 2014.
- 355 [27] Jeffrey Curtis Schlimmer. *Concept acquisition through representational adjustment*. PhD thesis,
356 University of California, Irvine, 1987.
- 357 [28] Yufeng Shang, Tao Liu, Yongchang Wei, Jingfeng Li, Liang Shao, Minghui Liu, Yongxi Zhang,
358 Zhigang Zhao, Haibo Xu, Zhiyong Peng, et al. Scoring systems for predicting mortality for
359 severe patients with COVID-19. *EClinicalMedicine*, 24:100426, 2020.
- 360 [29] A. J. Six, B. E. Backus, and J. C. Kelder. Chest pain in the emergency room: value of the heart
361 score. *Netherlands Heart Journal*, 16(6):191–196, 2008.
- 362 [30] Nataliya Sokolovska, Yann Chevalere, Karine Clément, and Jean-Daniel Zucker. The fused
363 lasso penalty for learning interpretable medical scoring systems. In *2017 International Joint
364 Conference on Neural Networks (IJCNN)*, pages 4504–4511. IEEE, 2017.
- 365 [31] Nataliya Sokolovska, Yann Chevalere, and Jean-Daniel Zucker. A provable algorithm for
366 learning interpretable scoring systems. In *International Conference on Artificial Intelligence
367 and Statistics*, pages 566–574. PMLR, 2018.
- 368 [32] Martin Than, Dylan Flaws, Sharon Sanders, Jenny Doust, Paul Glasziou, Jeffery Kline, Sally
369 Aldous, Richard Troughton, Christopher Reid, and William A Parsonage. Development and
370 validation of the emergency department assessment of chest pain score and 2h accelerated
371 diagnostic protocol. *Emergency Medicine Australasia*, 26(1):34–44, 2014.
- 372 [33] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal
373 Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- 374 [34] Berk Ustun, Lenard A Adler, Cynthia Rudin, Stephen V Faraone, Thomas J Spencer, Patricia
375 Berglund, Michael J Gruber, and Ronald C Kessler. The world health organization adult
376 attention-deficit / hyperactivity disorder self-report screening scale for dsm-5. *JAMA Psychiatry*,
377 74(5):520–526, 2017.
- 378 [35] Berk Ustun and Cynthia Rudin. Supersparse linear integer models for optimized medical scoring
379 systems. *Machine Learning*, pages 1–43, 2015.
- 380 [36] Berk Ustun and Cynthia Rudin. Optimized risk scores. In *Proceedings of the 23rd ACM
381 SIGKDD international conference on knowledge discovery and data mining*, pages 1125–1134,
382 2017.
- 383 [37] Berk Ustun and Cynthia Rudin. Learning optimized risk scores. *J. Mach. Learn. Res.*, 20:150–1,
384 2019.
- 385 [38] Berk Ustun, Stefano Traca, and Cynthia Rudin. Supersparse linear integer models for predictive
386 scoring systems. In *Workshops at the Twenty-Seventh AAAI Conference on Artificial Intelligence*,
387 2013.
- 388 [39] Caroline Wang, Bin Han, Bhrij Patel, and Cynthia Rudin. In Pursuit of Interpretable, Fair
389 and Accurate Machine Learning for Criminal Recidivism Prediction. *Journal of Quantitative
390 Criminology*, pages 1–63, 2022.
- 391 [40] Piotr Wasilewski, Bartosz Mruk, Samuel Mazur, Gabriela Póttorak-Szymczak, Katarzyna
392 Sklinda, and Jerzy Walecki. COVID-19 severity scoring systems in radiological imaging—a
393 review. *Polish journal of radiology*, 85(1):361–368, 2020.
- 394 [41] Feng Xie, Bibhas Chakraborty, Marcus Eng Hock Ong, Benjamin Alan Goldstein, Nan Liu, et al.
395 Autoscore: A machine learning–based automatic clinical score generator and its application to
396 mortality prediction using electronic health records. *JMIR medical informatics*, 8(10):e21798,
397 2020.
- 398 [42] Jiaming Zeng, Berk Ustun, and Cynthia Rudin. Interpretable classification models for recidivism
399 prediction. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 180(3):689–
400 722, 2017.
- 401 [43] Chi Zhang, Ling Qin, Kang Li, Qi Wang, Yan Zhao, Bin Xu, Lianchun Liang, Yanchao Dai,
402 Yingmei Feng, Jianping Sun, et al. A novel scoring system for prediction of disease severity in
403 COVID-19. *Frontiers in cellular and infection microbiology*, 10:318, 2020.

404 [44] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of*
405 *the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.

406 **Checklist**

- 407 1. For all authors...
- 408 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
409 contributions and scope? [Yes]
- 410 (b) Did you describe the limitations of your work? [Yes] See the Limitations section at the
411 end of Section 2
- 412 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See
413 Appendix G6. Even if a model is interpretable, it can still have negative societal bias
414 (though it is easier to check for such biases with scoring systems), and looking at a
415 variety of models from the pool could help find models that are more fair.
- 416 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
417 them? [Yes]
- 418 2. If you are including theoretical results...
- 419 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 420 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix C
- 421 3. If you ran experiments...
- 422 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
423 mental results (either in the supplemental material or as a URL)? [Yes] The code and
424 README are included as part of the supplemental material. Data links are included in
425 the Appendix D.1
- 426 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
427 were chosen)? [Yes] We perform 5-fold CV as specified in Section 4. Hyperparameters
428 are already specified (default values) in Algorithm 1 of Section 3.
- 429 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
430 ments multiple times)? [Yes] See Section 4.1. Error bars are included for the 5-fold
431 CV.
- 432 (d) Did you include the total amount of compute and the type of resources used (e.g., type
433 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix D.2
- 434 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 435 (a) If your work uses existing assets, did you cite the creators? [Yes] See Appendix D.3
436 and the References
- 437 (b) Did you mention the license of the assets? [Yes] See Appendix D.3
- 438 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
439 Our code is included as part of the supplementary material.
- 440 (d) Did you discuss whether and how consent was obtained from people whose data you're
441 using/curating? [N/A]
- 442 (e) Did you discuss whether the data you are using/curating contains personally identifiable
443 information or offensive content? [N/A]
- 444 5. If you used crowdsourcing or conducted research with human subjects...
- 445 (a) Did you include the full text of instructions given to participants and screenshots, if
446 applicable? [N/A]
- 447 (b) Did you describe any potential participant risks, with links to Institutional Review
448 Board (IRB) approvals, if applicable? [N/A]
- 449 (c) Did you include the estimated hourly wage paid to participants and the total amount
450 spent on participant compensation? [N/A]