

# Informed Injury Prediction in Elite Football: Decision Theory meets Machine Learning

Manuel Huth<sup>1,2</sup>, Berta Canal-Simón<sup>3,4</sup>, Eva Ferrer<sup>5,6</sup>, Gil Rodas<sup>5,6,7</sup>,  
Xavier Yanguas<sup>5</sup>, Jan Hasenauer<sup>1,2†</sup>, and Juan R González<sup>3,8,9†\*</sup>

<sup>1</sup> Life and Medical Sciences (LIMES) Institute, University of Bonn, Bonn, Germany

<sup>2</sup> Bonn Center for Mathematical Life Sciences, University of Bonn, Bonn, Germany

<sup>3</sup> Barcelona Institute for Global Health (ISGlobal), Barcelona, Spain.

<sup>4</sup> Made of Genes, Barcelona, Spain.

<sup>5</sup> Medical Department of Football Club Barcelona (FIFA Medical Centre of Excellence), and Barça Innovation Hub of Football Club Barcelona, Barcelona, Spain.

<sup>6</sup> Sports and Exercise Medicine Unit, Hospital Clinic and Sant Joan de Déu, Barcelona, Spain.

<sup>7</sup> Leitat technological Center, Terrassa, Spain.

<sup>8</sup> CIBER in Epidemiology and Public Health (CIBERESP), Barcelona, Spain.

<sup>9</sup> Department of Mathematics, Universitat Autònoma de Barcelona (UAB), Barcelona, Spain.

† Authors contributed equally; \* Correspondence: [juan.gonzalez@isglobal.org](mailto:juan.gonzalez@isglobal.org)

## 1 Abstract

2 Injuries in elite sports disrupt team performance, shorten careers, and incur significant finan-  
3 cial costs, highlighting the critical need for accurate predictions to inform optimal decisions  
4 that effectively prevent injuries. Existing approaches to injury prediction fail to account  
5 for cumulative risk, overlook injury severity, lack reliable probability calibration, and omit  
6 statistically guided decision thresholds. Here, we present a novel injury prediction frame-  
7 work integrating risk accumulation via survival analysis with machine learning, probability  
8 beta calibration, and statistical decision theory. Using a unique dataset spanning four sea-  
9 sons from FC Barcelona's women's team, we demonstrate that our framework outperforms  
10 standard classifiers, yielding superior discrimination ability. Our framework identifies fatigue-  
11 related measures as key injury predictors and incorporates flexible thresholds based on match  
12 importance and decision-maker certainty, improving player availability. Scalable and trans-  
13 ferable to other sports, this framework bridges academic research and practical deployment,  
14 empowering sports organizations to optimize player performance and long-term outcomes.

## 15 Introduction

16 Injuries in elite sports pose a significant challenge with wide-ranging implications, including  
17 disruptions to club performance, shortened player careers, and destabilized long-term finan-  
18 cial sustainability<sup>1-4</sup>. During the 2023/24 season alone, 4,123 injuries were recorded across  
19 the top five European leagues (England, Spain, Italy, Germany, and France), resulting in a  
20 financial burden of around 732.02 million euros<sup>5</sup>.

21 Efforts to predict daily injury risk using player tracking data and machine learning have  
22 emerged as a key strategy for mitigating these impacts<sup>6-9</sup>. These approaches leverage pre-  
23 dictive features, such as training loads, sleep quality, demographics, or psychological factors  
24 to anticipate injuries. Yet, significant challenges remain: a recent review by Bullock et al.<sup>10</sup>  
25 found that most existing models lack adequate performance assessment, probability cali-  
26 bration, and transparency, rendering them unsuitable for practical deployment. Similarly,  
27 Leckey et al.<sup>11</sup> highlight the misalignment between methodological development and practi-  
28 cal applicability, emphasizing that research teams must integrate both technical and domain  
29 expertise as even robust machine learning models are ineffective if they fail to address the  
30 real needs of practitioners.

31 We identify three additional gaps that hinder current injury prediction frameworks. First, ex-  
32 isting machine learning classifier-based approaches<sup>6-9;12;13</sup> overlook the cumulative nature of  
33 injury risk. Specifically, these methods fail to capture how injury probability increases mono-  
34 tonically with the duration of exposure during training or matches. For example, a player  
35 participating for 70 minutes should inherently face a lower injury risk than one playing for 90  
36 minutes. Second, models lack day-specific tailored decision thresholds that reflect the varying  
37 importance of matches versus training sessions. For example, a player with moderate injury  
38 risk might strategically rest during a training day but participate in a match due to its higher  
39 strategic and financial value. Constructing these thresholds requires reliable probability cali-  
40 bration—that is, ensuring that predicted probabilities align with observed outcomes. For  
41 example, if a model predicts a 5% chance of injury, then approximately 5% of those cases  
42 should indeed result in an injury<sup>14</sup>. Unfortunately, traditional machine learning models of-  
43 ten fail to achieve this level of calibration without further adjustments<sup>15</sup>. Consequently, any  
44 robust framework must include mechanisms to correct miscalibrated probabilities. Third,  
45 commonly used metrics like Recall<sup>6-8</sup> or the F1 score<sup>6;9</sup> fail to account for the severity of in-  
46 juries, measured in injury duration. While detecting a single long-term injury may be highly  
47 beneficial for team performance and availability, traditional evaluation metrics penalize false  
48 positives without incorporating the practical impact of injury severity.

49 To address these limitations, we introduce a novel injury prediction framework that inte-  
50 grates survival analysis with machine learning, probability calibration, and practical decision-  
51 making strategies derived from statistical decision theory - all designed to meet the real-world

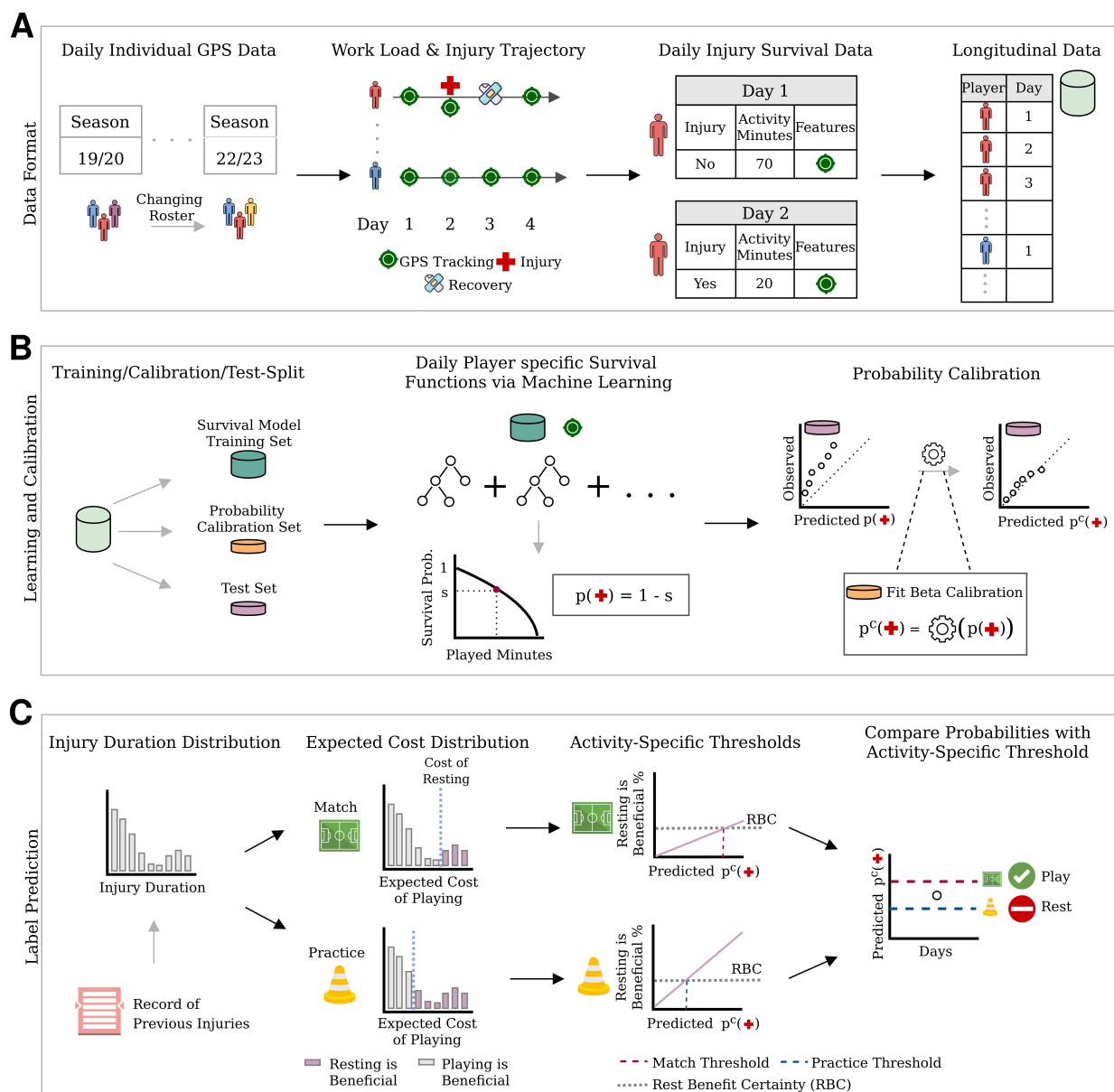
52 needs of practitioners. Specifically, our approach models injury risk by treating each activ-  
53 ity as a survival observation, enabling the estimation of injury probabilities over time using  
54 methods from survival analysis that naturally account for injury risk accumulation over  
55 time<sup>16–18</sup>. We combine GPS-derived training load features to create a comprehensive and  
56 individualized assessment of injury risk. To ensure that predicted probabilities are reliable  
57 for decision-making, we apply beta calibration<sup>19</sup>, a technique widely used in fields such as  
58 oncology research where accurate probabilities are essential<sup>20;21</sup>. Furthermore, we derive ac-  
59 tionable decisions by incorporating day-specific valuations for matches and training sessions,  
60 alongside a pre-specified level of rest benefit certainty (RBC) that reflects decision-makers’  
61 required confidence to rest a player.

62 We evaluate our framework using a dataset spanning four seasons (2019/20 to 2022/23) from  
63 Fútbol Club Barcelona’s (FC Barcelona, Spain) first women’s team<sup>22;23</sup>. Injuries are consid-  
64 ered to reflect a team’s time loss of player availability. Injury evaluations were conducted by  
65 the team’s medical physician in collaboration with the FC Barcelona medical department,  
66 adhering to standardized diagnosis and return-to-play protocols established in the club’s  
67 guidelines<sup>24;25</sup>. The analysis specifically targeted non-contact injuries involving muscles, ten-  
68 dons, ligaments, and cartilage. First, we assess the robustness of survival-based models  
69 compared to classical machine learning classifiers<sup>26–28</sup> over 400 training, calibration, and test  
70 splits from the 2019/20, 2020/21, and 2021/22 seasons. Model performance is evaluated using  
71 a player availability gain metric, which accounts for injury durations to provide a practical  
72 assessment of prediction impact. Second, we identify the most important predictive features  
73 across these splits to better understand key contributors to injury risk. Finally, we validate  
74 the best-performing model and selected features on the previously unseen 2022/23 season  
75 using a rolling prediction approach, simulating real-world conditions for decision-making.

## 76 Results

### 77 Integrating Risk Accumulation, Calibration, Match Valuations and 78 Decision Certainty in Injury Prediction

79 An effective injury prediction framework for clubs, particularly for coaches and medical staff  
80 to which we refer as team decision-makers, must meet three criteria: (1) capture the non-  
81 linear accumulation of injury risk over training time, (2) provide accurate probability scores  
82 reflecting true injury risk, and (3) offer interpretable thresholds to translate probabilities into  
83 actionable decisions, such as predicting to rest a player or to let her play.



**Figure 1: Data Format, Learning and Label Prediction.** (A) Data Format: Daily GPS-based training load and injury records were collected over 4 seasons with a changing player roster. Each day is treated as an individual survival observation, where the injury status serves as the censoring indicator and activity minutes as survival time. Combining daily observations per player creates a longitudinal survival dataset. (B) Learning and Calibration: The data is split into training, calibration, and test sets. The training set is used to fit the survival model, which outputs the probability of injury given planned activity time. The calibration set is used to refine injury probabilities via beta calibration using observed injury indicators. (C) Label Prediction: Injury distributions are estimated from past records. For match and practice days, expected costs of resting and playing are compared for each historical injury duration. If the Rest Benefit Certainty (RBC), e.g., 50%, favors resting, the player is predicted to rest; otherwise, they play.

84 In this study, we introduce an injury prediction framework that meets these criteria by incor-  
85 porating risk accumulation over time and dynamically adjusting decision thresholds based  
86 on match and practice session importance as well as required certainties of decisions: (1) We  
87 interpret the machine learning task as a daily survival analysis for each player (Figure 1A).  
88 Survival in this context refers to the time until an injury in minutes. Each activity day  
89 represents a survival observation, with played minutes treated as the time indicator and the  
90 injury label acting as the censoring indicator. Unlike traditional methods that treat this as  
91 a pure classification task<sup>6;13;29</sup>, survival analysis captures the dynamics of injury risk accu-  
92 mulation over the activity time by a monotonic decreasing survival curve. As main survival  
93 model, we use the Extreme Gradient Boosting Survival Embedding (XGBSE)<sup>16</sup>. For further  
94 comparison, we also include the Accelerated Failure Time<sup>17</sup>, Coxnet<sup>18</sup>, and boosted Cox<sup>30</sup>  
95 survival models. The machine learning survival model is trained using the input features, e.g.  
96 accumulated GPS data, a match indicator and age, on the training set. This model generates  
97 survival curves for unseen observations based on these input features (Figure 1B). Survival  
98 curves evaluated at the planned training duration yield an estimate of the probability of  
99 sustaining no injury up to the planned training duration, which is then used to compute  
100 the probability of injury given the planned training duration. By construction of the mono-  
101 tone survival curves, longer training durations naturally lead to higher injury risk, reflecting  
102 risk accumulation over time. In contrast, standard classification models like standard Light-  
103 GBM<sup>26</sup>, which treat played minutes as an ordinary feature, fail to reflect this cumulative  
104 risk. For LightGBM, a monotonization feature has been introduced that could, in principle,  
105 deal with the monotone risk accumulation. We benchmark the monotone feature and other  
106 standard machine learning classifier<sup>26–28</sup> against the survival analysis approach within our  
107 results. For model training, calibration and evaluation, the data is split into a training, a  
108 calibration, and a test set.

109 (2) To align predicted probability scores with the actual injury risk, we calibrate predicted  
110 probabilities using a separate calibration set and a beta calibration model<sup>19</sup>. This calibra-  
111 tion ensures that the reported probabilities reflect actual injury risk and thus are actionable  
112 for team decision-makers (Figure 1B). Without calibration, thresholds for actionable deci-  
113 sions may become unreliable, leading to overly conservative or risky choices. For instance,  
114 miscalibrated probabilities might underestimate injury risk, exposing players to harm, or  
115 overestimate it, sidelining key players unnecessarily. Beta calibration is particularly suitable  
116 due to its robustness to skewed probability distributions<sup>19</sup>. We validate the calibration for our  
117 data by comparing the calibrated injury probabilities against the uncalibrated probabilities.

118 (3) Calibrated probabilities provide a robust foundation for playing-or-resting decision-making  
119 but must be paired with actionable decision thresholds. Our framework leverages pre-  
120 specified activity valuations that can be tailored to individual players and specific days.  
121 These valuations reflect the importance of participation on any given day, enabling decisions

122 that align with team priorities. A higher valuation is reflected in a higher probability thresh-  
123 old required to rest a player. For instance, these valuations can reflect the higher stakes of  
124 critical matches or lower priority of routine training (Figure 1C).

125 Furthermore, coaches and medical staff can specify a minimum certainty level required to rest  
126 a player - reflecting their risk preferences. This is achieved by calculating the expected costs  
127 of playing - which depend on the distribution of potential injury durations - and comparing  
128 them to the fixed cost of resting. This comparison quantifies the likelihood that resting a  
129 player will result in a player Availability Gain (AG). For example, based on predicted injury  
130 risk and observed injury durations, decision-makers might agree on resting a player if it is  
131 beneficial in 50% of injury duration scenarios. Decision-makers can use this level of certainty  
132 to guide their actions, which we label as Rest Benefit Certainty (RBC).

133 To introduce these match valuation and RBC concepts further to the reader, we have included  
134 examples with concrete numbers in the Supplementary Material.

## 135 **Collection, Assembly and Curation of Training Data**

136 Our study compiled a dataset for FC Barcelona's first women's team over four consecutive sea-  
137 sons (2019/2020 to 2022/2023). This integrated dataset combines daily training load records,  
138 match schedule information, injury reports, and international duty data. When training data  
139 were missing due to players' international duty, match participation (e.g., minutes played)  
140 was used for imputation of training load. Daily injury records - detailing occurrence and  
141 duration - further enriched our dataset, allowing us to examine links between training, per-  
142 formance, and injury risk. All data were harmonized to a common timeline, providing a  
143 robust foundation for our analysis.

144 The dataset includes 34 players, with turnover as some players left and new players joined the  
145 club during this period (Figure 2A). This turnover introduces variability due to unobserved  
146 player differences, which could affect model performance and requires results to be interpreted  
147 accordingly. The players' ages ranged from 17 to 33 years, with a slight year-to-year increase  
148 due to natural aging (Figure 2B). This increase in ages was offset by younger players joining  
149 the team in 2022/23, resulting in marginal changes to the age distribution. The frequency of  
150 matches and practices varied over time due to injuries, international duties, personal breaks,  
151 and the COVID-19 disruption during the 2019/20 season. While practices outnumbered  
152 matches (Figure 2C), the distribution of covered kilometers per day was stable across seasons  
153 and higher for matches compared to practices (Figure 2D). These factors were incorporated  
154 into the survival model using a match indicator and aggregated training load variables from  
155 days prior to the day of activity. Over four seasons, 83 non-contact injuries were recorded  
156 (0.62% of all data points), with a pronounced spike in 2021/2022 (31 injuries) compared  
157 to the prior season (14 injuries; Figure 2E). It is striking that injuries occur at a high rate

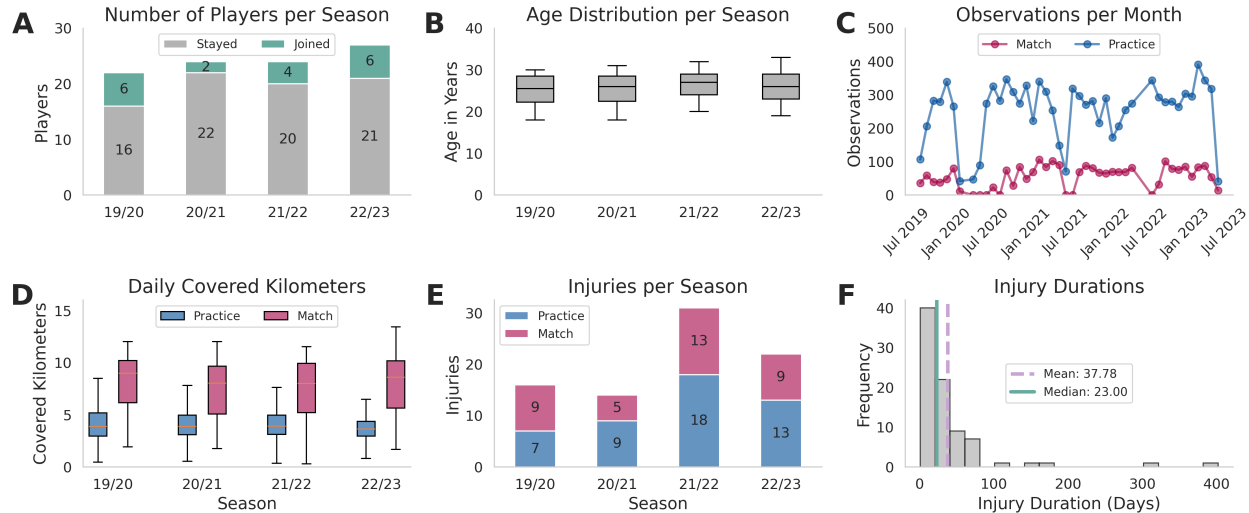


Figure 2: **Summary of player characteristics, observations, and injuries across seasons.** (A) Number of players per season, with returning players shown in gray and new players in green. (B) Boxplot of age distribution of players per season. (C) Monthly number of observations, split into matches (red) and practices (blue). (D) Boxplot of daily covered kilometers, shown as boxplots. (E) Number of injuries per season, separated into match injuries (red) and practice injuries (blue). (F) Distribution of missed days per injury.

158 during matches, even though the total practice time per season is approximately four times  
 159 greater than the match time (Supplementary Figure S1) suggesting an increased non-contact  
 160 injury risk during matches. The high class imbalance and the shift in injury label distribution  
 161 present significant challenges for injury prediction models as these fluctuations can undermine  
 162 the stability and predictive performance of models trained on historical data. Injury durations  
 163 ranged widely, with a mean of 37.8 days and a median of 23 days (Figure 2F). This disparity  
 164 highlights the influence of long-term injuries, such as ligament injuries, and importance of  
 165 analyzing the full injury duration distribution, rather than relying on averages, for informed  
 166 decision-making.

## 167 **Survival Models Outperform Standard ML Classifiers in Discrimi-** 168 **nation Ability**

169 To evaluate the discrimination ability of our proposed framework, we compared its perfor-  
 170 mance against alternative survival models and standard machine learning classifiers. Here,  
 171 discrimination refers to a model’s capacity to distinguish between positive (injury) and neg-  
 172 ative (no injury) cases independently of any decision threshold. This approach enables us to  
 173 assess model performance without the confounding effects of threshold selection. For compar-  
 174 ison, we have chosen LightGBM<sup>26</sup> as tree-based Extreme Gradient Boosting has performed

175 best in previous injury prediction attempts<sup>7-9</sup>, and additionally k-nearest-Neighbors<sup>27</sup> and  
176 Logistic Regression<sup>28</sup> for a broader scope of the comparison. We also incorporate two addi-  
177 tional variants of LightGBM. The first enforces a monotonic relationship for playing time so  
178 that higher training loads correspond to increased injury risk. The second leverages Light-  
179 GBM as a tree embedding, which is then fed into a Logistic Regression model—mimicking  
180 the approach behind XGBSE for a standard machine learning classifier.

181 We quantified discrimination using the Mean Average Precision (MAP) and the Area Under  
182 the Receiver Operating Characteristic Curve (AUROC). While AUROC has been widely  
183 employed in injury prediction studies<sup>6-8</sup>, Mean Average Precision (MAP)—which jointly  
184 accounts for both precision and recall—has received relatively little attention. This is despite  
185 its proven effectiveness in handling highly imbalanced datasets<sup>31;32</sup>, a common challenge in  
186 injury prediction where healthy player records far exceed those of injuries. As a baseline,  
187 we employed a prior mean classifier that leverages the inherent class frequencies to ensure  
188 meaningful comparisons. Data were partitioned using a month-wise block splitting strategy  
189 (Figure 3A) to preserve temporal dependencies. This approach split the data into distinct  
190 training, calibration, and test sets, and the splitting process was repeated 400 times to ensure  
191 robust performance estimates.

192 XGBSE achieved a median MAP of 0.071 - an approximately 11-fold improvement over the  
193 baseline mean classifier (Figure 3B) - and a median AUROC of 0.761, improving the baseline  
194 mean classifier substantially.

195 In further comparisons, XGBSE outperformed the best standard classifier, LightGBM (which  
196 achieved a median MAP of 0.021), by a factor of roughly 3.44. Although the Accelerated  
197 Failure Time model demonstrated the second-best MAP (median: 0.048) and the highest  
198 AUROC (median: 0.780), its performance was closely followed by XGBSE in terms of AU-  
199 ROC. The similarity in performance between these models is unsurprising given that XGBSE  
200 leverages the Accelerated Failure Time model as a feature transformer.

201 Overall, our results show for both considered metrics that survival models, particularly XG-  
202 BSE, offer superior discrimination ability compared to standard machine learning classifiers.

## 203 **Beta Calibration Aligns Predicted Probabilities with True Injury** 204 **Risk**

205 To assess if the models provide interpretable probability estimates that are well-aligned with  
206 actual injury risk, we analyzed how the predicted probabilities generated by XGBSE aligned  
207 with the observed proportions of injuries. By performing calibration on an independent  
208 calibration set - distinct from both the training and test sets - our approach ensures an  
209 unbiased alignment of predicted probabilities with the true injury risk for each of the 400

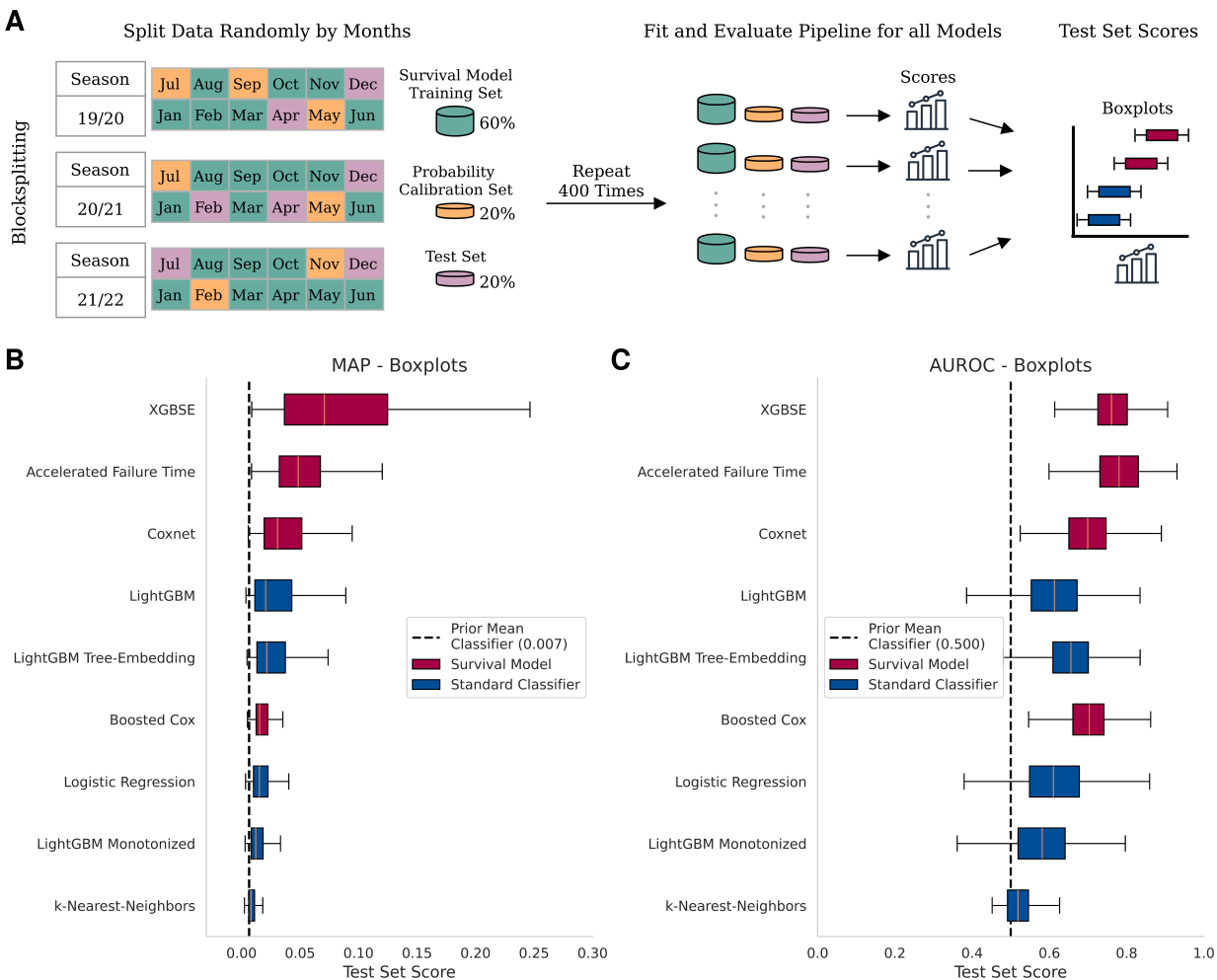


Figure 3: **Discrimination Performance.** Survival-based models are shown in red, and standard machine learning classifiers in blue. (A) The sample is split into training, calibration, and test sets, repeated 400 times with different seeds using a month-wise block-splitting such that observations from the same month stay within the same set. For each split, the framework is fitted, and the metrics are computed. (B, C) Boxplots of Mean Average Precision (MAP) and Area Under the Receiver Operating Characteristic Curve (AUROC) across 400 splits.

210 splits (Figure 4A). We computed average calibration curves by deriving a calibration curve  
 211 for each split and then taking their point-wise average.

212 Our analysis showed that the uncalibrated XGBSE model systematically underestimated the  
 213 observed proportions of injuries, as evidenced by the average calibration curve (Figure 4B).  
 214 After calibration, the average calibration curve closely aligned with the observed injury rates,  
 215 demonstrating improved accuracy and reliability in the predicted probabilities.

216 We observed similar behavior with other methods, such as the Accelerated Failure Time

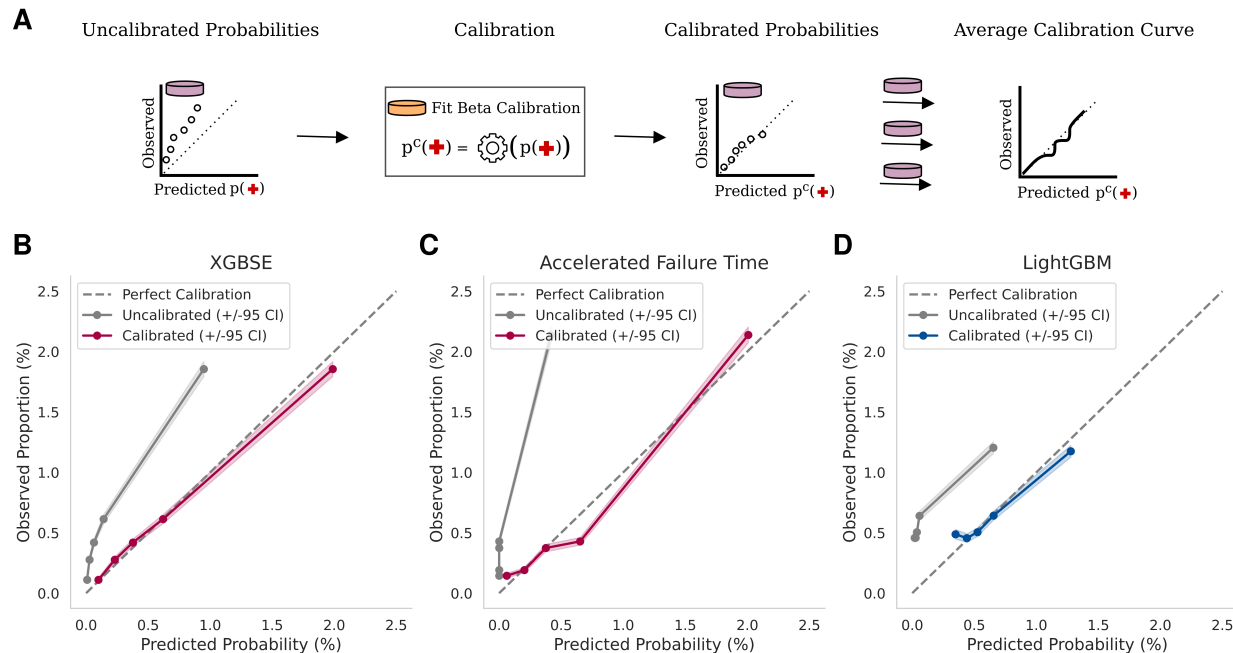


Figure 4: **Calibration Performance.** Average calibration curves for uncalibrated and calibrated predictions, using 5 quantile bins per curve. Averages are computed point-wise across 400 splits, with the 45° line indicating perfect calibration. **(A)** Representation of calibration workflow. **(B)** For the XGBSE model. **(C)** For the Accelerated Failure Time model. **(D)** For LightGBM.

217 survival model (Figure 4C) and LightGBM (Figure 4D), demonstrating that calibration gen-  
 218 erally improves risk alignment. However, these results, combined with our earlier findings  
 219 on discrimination ability, indicate that while beta calibration can effectively align risk es-  
 220 timates, a robust underlying model is essential for achieving strong predictive power. For  
 221 the survival models, calibration involved only monotone transformations that did not alter  
 222 the MAP or AUROC while for LightGBM, only some transformations were non-monotone  
 223 (Supplementary Figure S2) such that results stay qualitatively the same.

## 224 **Uncertainty-Aware Cut-offs Improve Player Availability Across Match** 225 **Valuations**

226 After examining the discrimination and calibration capabilities of our proposed framework,  
 227 we evaluated its performance in decision-making scenarios to test its applicability as decision  
 228 support for coaches and medical staff (Figure 5A). Specifically, we analyzed three match  
 229 valuation scenarios:  $R = 2, 5, 10$ . Here,  $R$  represents the relative match valuation, quantifying  
 230 how much more important a match is compared to a training session. The calculations

231 were conducted for five RBC thresholds ranging from 50–90%, representing varying levels  
232 of confidence required to rest a player. The primary evaluation metric, Mean Availability  
233 Gain (MAG), quantifies the trade-off between the costs of unnecessary rest and the benefits  
234 of detecting injuries that lead to lost match and practice days. A more detailed explanation  
235 is provided in the Methods section.

236 Each of the split results were multiplied by a scaling factor to match the number of obser-  
237 vations from the 2022/2023 season making the results comparable with the seasonal analysis  
238 of 2023. Across all match valuations and RBC thresholds, we observed consistently positive  
239 MAGs (Figure 5B).

240 Specifically, a RBC range of 50–70% consistently yielded the highest MAGs across all match  
241 valuations. For higher match valuations, the differences in MAGs across RBCs were smaller,  
242 suggesting greater flexibility in threshold selection. Notably, MAGs within the 50–70% RBC  
243 range were not statistically different from each other at the 95% confidence level across all  
244 match valuations.

245 We analyzed the precision, recall and the F1-score, to evaluate the trade-offs in predicting  
246 injuries within a highly imbalanced dataset while ignoring the impact of injury durations  
247 (Supplementary Figure S3). For a match valuation of 5 and a RBC of 60%, precision indicates  
248 that around 1 in 11 (9.1%) predicted rest days corresponded to actual injuries, while recall  
249 shows the model detected around 1 in 12 (8.1%) injury events.

## 250 **Cumulative Load Emerges as the Strongest Predictor**

251 To better understand the mechanics of our method’s predictive performance and to assess  
252 feature relevance, we examined the influence of individual features on predictions (Figure 6).  
253 Specifically, we calculated two metrics for each sample split: the average gain in the loss func-  
254 tion and the frequency with which each feature was selected for splits across all trees in the  
255 ensemble. The medians across all sample splits were used to summarize results (Figure 6A).

256 The strongest predictor for both metrics was the exponentially weighted moving average of  
257 distance covered (in meters) over the past 21 days, capturing the cumulative effects of player  
258 load from prior playing activities. Following this, high-speed running intensity (HSR per  
259 minute) and the frequency of accelerations and decelerations on the day of activity emerged  
260 as the next most influential factors. In contrast, the cumulative acceleration and deceleration  
261 workload over the past 21 days had minimal impact, suggesting that while acute accelerations  
262 and decelerations on the day of activity contribute to injury risk, their accumulated effect  
263 over time plays a comparatively minor role in fatigue-induced injury susceptibility.

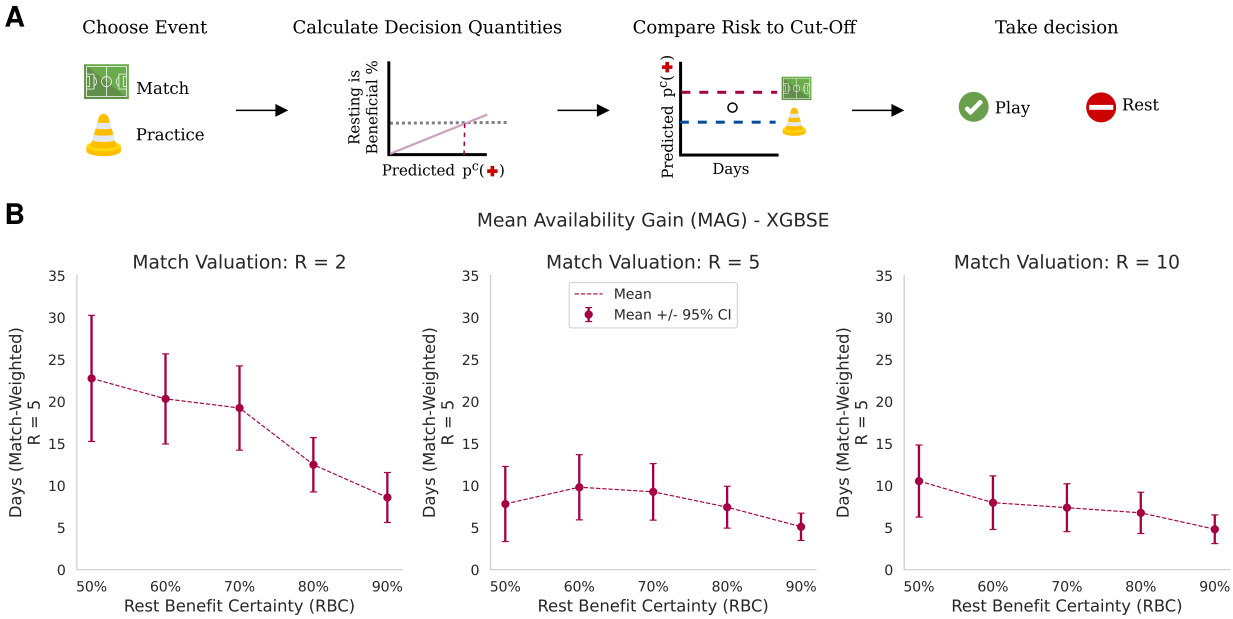


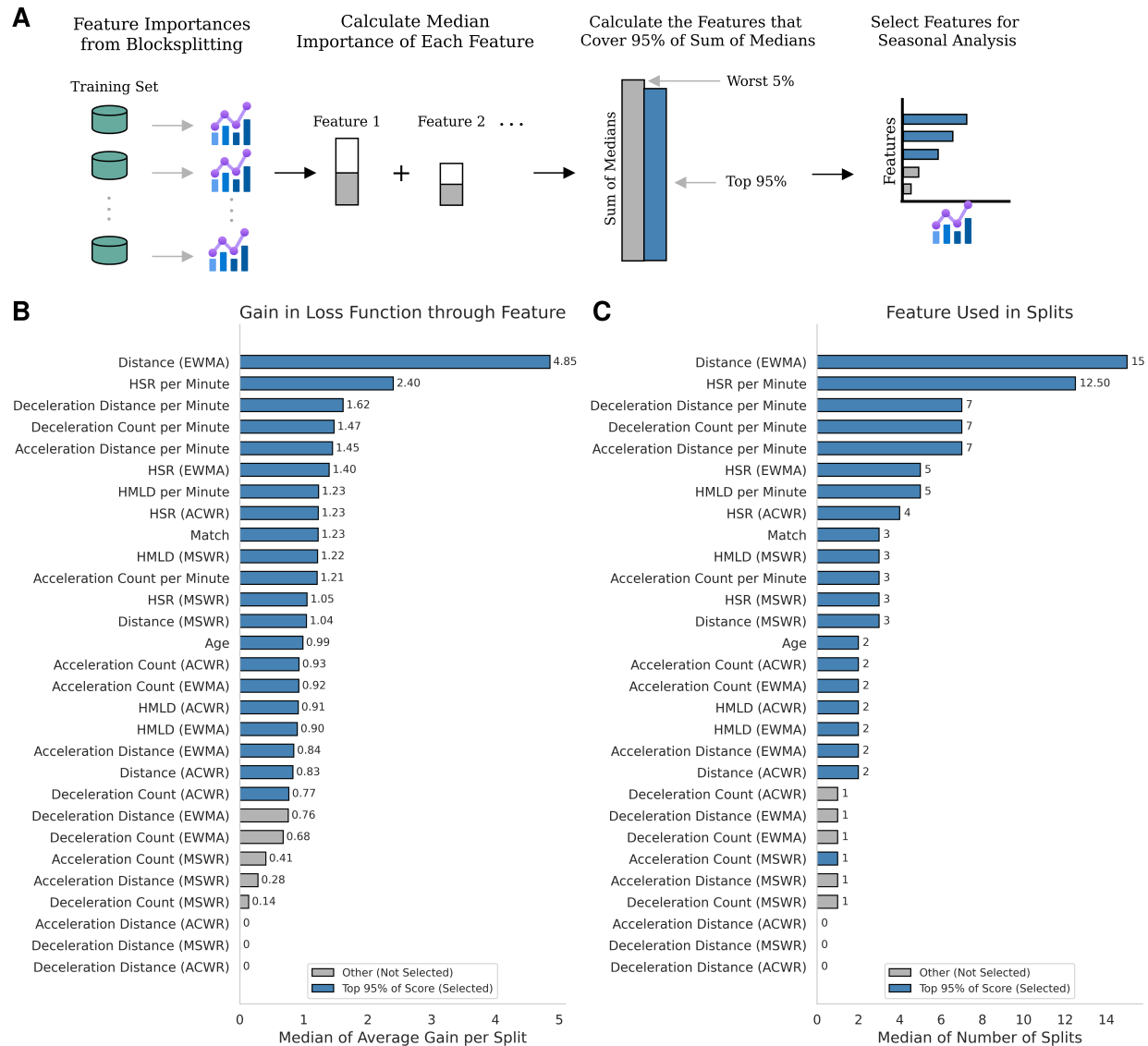
Figure 5: **Rest Benefit Analysis of Player Availability.** (A) Representation of Decision Process. (B) Mean Availability Gain as a function of Rest Benefit Certainty (RBC) and match valuation, calculated over 400 samples for the XGBSE model. Confidence intervals represent the 95% confidence range for the mean, obtained via a normal approximation.

## 264 Seasonal Prediction Validates Framework Performance with Posi- 265 tive Availability Gains

266 To evaluate the real-world applicability of our proposed framework in daily coaching scenar-  
267 ios, we applied the framework using the XGBSE model to predict player outcomes for the  
268 previously unseen 2022/2023 season. Predictions were generated in a weekly rolling fash-  
269 ion to simulate a realistic scenario where the model is re-trained weekly, incorporating the  
270 most recent data. This approach ensures that predictions are made using the latest trends  
271 and benefits from increasing training and calibration set sizes over time (Figure 7A). After  
272 obtaining all predictions, we aggregated the predicted labels into a test set for evaluating  
273 seasonal performance.

274 The features used for these predictions were selected based on the feature importance analysis  
275 of data from the 2019/2020 to 2021/2022 seasons (Figure 6A). To identify key predictors, we  
276 accumulated the medians for the gain in the loss function and selected features accounting  
277 for 95% of the cumulative importance.

278 To showcase how a seasonal trajectory for one player looks like, we show predictions for  
279 one player (Figure 7B) for a match valuation of  $R = 5$  and a RBC of 60%, that highlight  
280 the potential of the framework. In August 2022, the player suffered one injury, which was



**Figure 6: Feature Importance.** Feature importances were computed for the XGBSE model using the feature importances of the underlying Accelerated Failure Time model. Inter alia, Exponentially Weighted Moving Averages (EWMA), Acute Chronic Workload Ratio (ACWR), and Monotonic Workload Ratio (MWR) were used. **(A)** The average gain in the loss function when selecting a feature and the number of times a feature is used in a split were calculated. The median of both metrics was computed across the 400 splits. Features with the lowest medians that collectively contribute less than 95% of the total sum of medians for both metrics are identified as uninformative. **(B)** Median of the average gain per split. Features contributing to 95% of the total sum of medians are highlighted in blue. **(C)** Median of the number of times a feature is used in a split within the tree ensemble. Features contributing to 95% of the total sum of medians are highlighted in blue.

281 correctly predicted by the model. After recovering from the injury, the player participated  
282 in three matches that exceeded the practice threshold but not the match threshold without  
283 sustaining further injuries. Hence, the model correctly classified these events as non-rest  
284 days. In November 2022, the model classified another rest day that did not result in an  
285 injury; this prediction was close to the classification threshold. Overall, the player would  
286 have been rested three times during the season: once when an injury occurred, once two days  
287 before an injury, and once when no injury occurred but the predicted injury probability was  
288 close to the classification threshold.

289 The availability gain curve for the 2022/2023 season, starting at a RBC of 50%, followed a  
290 decreasing trend, consistent with the split-based results from the 2019/2020 to 2021/2022  
291 seasons (Figure 7C). Notably, RBCs of 50% and 60% yielded availability gains of around 20.1  
292 and 14.7 match-weighted days, respectively, while RBCs of 80% and 90% resulted in slightly  
293 negative availability gains of -2.4 and -1.2 match-weighted days. These findings suggest that  
294 higher RBC thresholds, which require higher probabilities to classify a player as requiring  
295 rest (Supplementary Figure S4), may fail to prevent major injuries.

296 To evaluate the detection of injuries and the rate of misclassifying non - injury days as injuries,  
297 we further analyzed factors contributing to availability gain. This analysis included unneces-  
298 sary rest days (false positives), saved days (prevented rest days due to detected injuries), as  
299 well as key performance metrics such as precision and recall (see Supplementary Figure S5).  
300 Across all match RBCs, the model consistently avoided predicting rest days during matches,  
301 demonstrating stable behavior. At a RBC threshold of 50%, the model identified 5 out of  
302 22 injuries (recall: 22.73%), while only 1 in 18 predicted rest days corresponded to an actual  
303 injury (precision: 5.44%).

304 Results for match valuations  $R = 2, 10$  were qualitatively the same and are reported within  
305 the Supplementary Figure S6 and Figure S7.

## 306 Seasonal Predictions Align with Observed Data

307 To evaluate how well the predicted probabilities align with observed injury data, we com-  
308 puted the probability mass function of the cumulative number of injuries for each day. Our  
309 results demonstrate that the observed data aligns well with the high-probability mass re-  
310 gions predicted by the model (Figure 7D). The predicted mean trajectory closely follows the  
311 trajectory of the observed number of injuries. This implies that the model is well calibrated.

312 We further assessed model performance using the AUROC metric, achieving a score of 0.74  
313 (Figure 7E). This value indicates a strong ability to discriminate between injury and non-  
314 injury events in the context of the high class imbalance with only 0.58% positive cases during  
315 the 2022/2023 season. Precision-recall analysis showed a Mean Average Precision (MAP) of

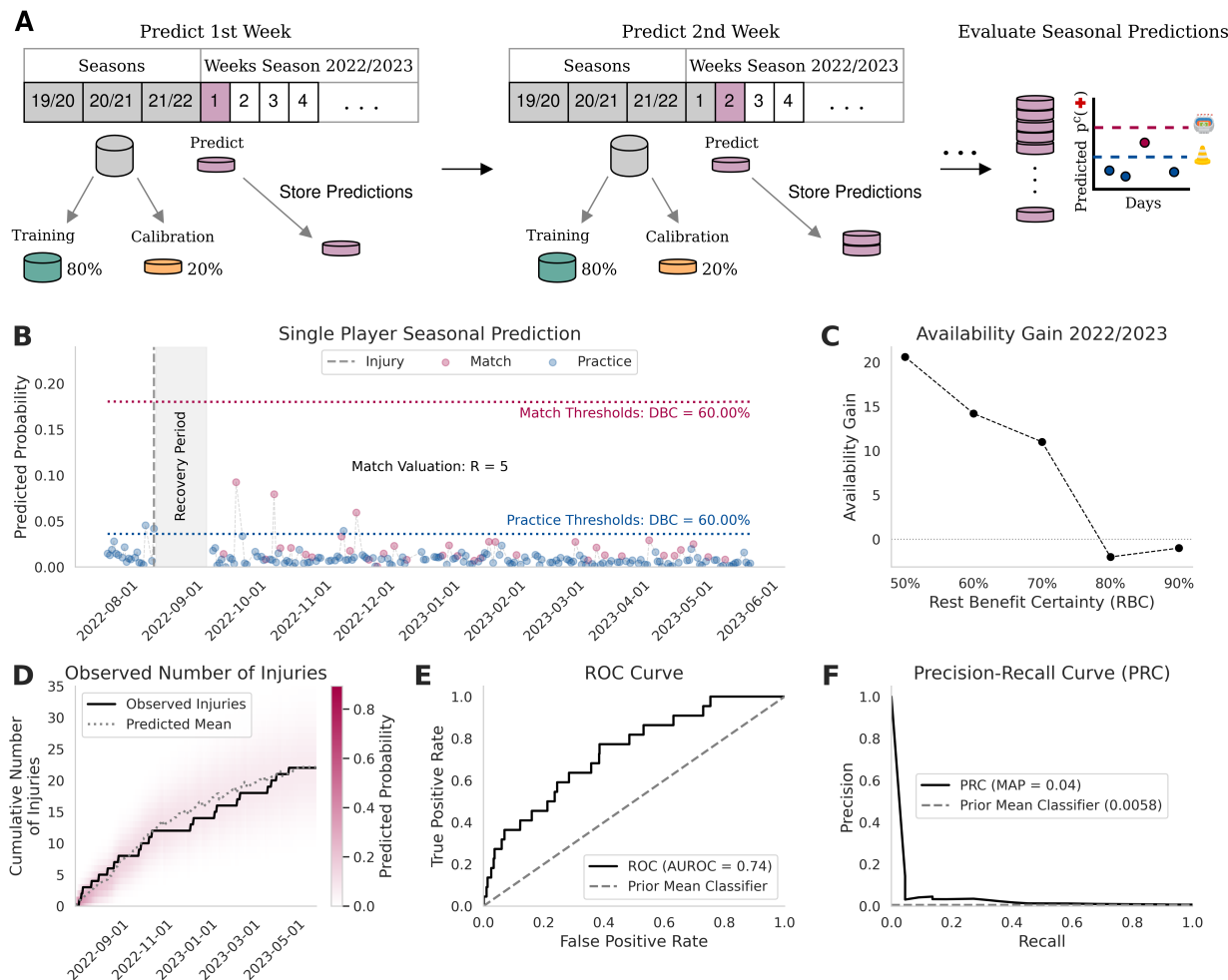


Figure 7: **Predicting the 2022/2023 Season.** (A) The previously unseen 2022/2023 season is predicted in a rolling weekly fashion. Seasons 2019–2022 are used as the baseline training set. Every 7 days, the model is re-trained using the baseline data and all available data up to that day, forming a training super set. Predictions for the subsequent week are made with the re-trained model. The training super set is randomly split into a training set and a calibration set, with 20% allocated to calibration. Weekly predictions are aggregated to evaluate the overall seasonal metrics. (B) Single-player results obtained with a match valuation of  $R = 5$  and a Rest Benefit Certainty (RBC) of 60%. (C) Availability gain as a function of differing RBC values, evaluated on the accumulated test set. (D) Predicted distribution of the cumulative number of injuries per day (red) with its mean compared to observed injuries. The distribution is computed via iterative convolution over days. (E, F) Receiver Operating Characteristic (ROC) curve and Precision-Recall (PR) curve evaluated on the test set.

316 0.04, representing an approximately 6.9-fold improvement over random guessing (Figure 7F).

## 317 Discussion

318 This study presented a novel injury prediction framework that integrates survival-based ma-  
319 chine learning, probability calibration, and actionable decision thresholds tailored to day-  
320 specific valuations. These actionable insights are derived using rest-or-play thresholds de-  
321 rived from statistical decision theory. Therefore, predicted probabilities must accurately  
322 reflect actual injury risks. Beta calibration significantly enhanced the reliability of predicted  
323 probabilities, ensuring alignment with true injury risks and enabling actionable decisions.

324 Using a unique dataset from FC Barcelona’s first women’s team, we demonstrated that our  
325 framework, leveraging survival-based approaches, significantly outperforms standard machine  
326 learning classifiers<sup>26–28</sup> in their ability to predict injury risk. The survival-based models can  
327 account for risk accumulation over time - a critical limitation of standard classifiers that are  
328 used in the literature<sup>6;13;29</sup>. Positive player availability gains - a metric, to the best of our  
329 knowledge, not previously reported in the literature - were observed across a range of match  
330 valuations and decision certainty thresholds, underscoring the framework’s robustness and  
331 flexibility in balancing injury prevention with team needs. The XGBSE model emerged as the  
332 top-performing method, validated through evaluation on three seasons of historical data and  
333 further applied to previously unseen data from the 2022/23 season, where it led to meaningful  
334 gains in player availability and closely aligned with the observed injury distribution.

335 We examined the most predictive features and identified fatigue-related measures, such as  
336 accumulated training distances of previous days, and same-day activity intensity emerged  
337 as the most influential predictors, showing that fatigue-related injury risk can be primarily  
338 predicted by previously long training sessions. Similarly, same-day intensity-related measures  
339 like acceleration and deceleration per minute account for a higher injury risk on that day.

340 The framework is highly flexible and, although demonstrated on professional soccer, is  
341 broadly applicable to other sports and domains. With the proliferation of wearable tracking  
342 devices not only among professional athletes but also amateurs and workers in physically  
343 demanding fields, our method offers a powerful tool for injury prediction in a variety of  
344 settings. The survival modeling approach seamlessly integrates features relevant to diverse  
345 injury mechanisms and domain-specific risk factors—ranging from biomechanical loads and  
346 psychological factors to specialized activity metrics. Even image or video data could straight-  
347 forwardly be incorporated to detect anomalies in movement patterns.

348 For other sport teams, we recommend re-training the model to capture differences in training  
349 regimes and player profiles. The computational efficiency of the XGBSE model - training on  
350 the entire dataset takes less than a second (Supplementary Figure S9) - makes re-training  
351 feasible even for large datasets.

352 Despite its strengths, our method has limitations. First, similar to standard evaluation

353 metrics like the F1 score or balanced accuracy, the availability gain metric does not represent a  
354 causal intervention effect. Specifically, (a) our evaluation does not fully account for scenarios  
355 where resting a player before an injury could have prevented it, as shown in the seasonal  
356 analysis, and (b) resting a player on a high-risk day may prevent an injury that could still  
357 occur later. Since (a) underestimates and (b) overestimates the real-world effect, the net  
358 impact remains unclear, precluding causal interpretation. Notably, while this limitation  
359 applies to all metrics, availability gain uniquely accounts for both injury occurrence and  
360 injury severity. Second, while our dataset from FC Barcelona's women's team spans four  
361 seasons, validation on data from other clubs would enhance generalizability. Differences  
362 in training, playing styles, and player demographics could provide further insights into the  
363 framework's robustness.

364 Future improvements could enhance both the methodological and practical aspects of the  
365 framework. Its fast computation enables real-time applications, allowing the integration  
366 of live match or training data for immediate insights and rapid, data-driven decisions -  
367 potentially even during matches or training. A long-term goal is the development of digital  
368 twins that simulate player states, team dynamics, and causal injury intervention effects in  
369 real time. By integrating our model within multiple modeling approaches - such as match  
370 simulations - these digital twins could optimize data-driven decision-making for teams.

371 Incorporating player-specific metrics, such as salary or market value loss, into the loss func-  
372 tion could help prioritize long-term injury prevention. Additionally, incorporating informa-  
373 tive features - such as heart rate data - and employing random effect models or Gaussian  
374 processes to capture temporal correlations in survival curves may further enhance personal-  
375 ized predictions.

376 Overall, this work demonstrates the capability and practicality of using survival-based ma-  
377 chine learning combined with probability calibration and statistical decision theory to predict  
378 injury risks and inform actionable decisions in elite sports. By incorporating day-specific  
379 valuations and decision-maker certainty thresholds, our approach bridges the gap between  
380 academic research and real-world deployment. Its scalability, flexibility, and transferabil-  
381 ity make it a valuable tool for sports organizations seeking to optimize player availability,  
382 performance, and long-term outcomes across diverse settings.

## 383 **Methods**

### 384 **Employed Models**

385 To predict injury risk in professional footballers, we employ both parametric and non-  
386 parametric survival models. Parametric models, such as the penalized Cox proportional

387 hazards model (Coxnet)<sup>18;33;34</sup>, specify a functional form for the hazard rate, offering well-  
388 understood statistical properties. Non-parametric models, including the Extreme Gradient  
389 Boosting Accelerated Failure Time model<sup>17;35</sup> and the boosted Cox model<sup>30;34</sup>, provide the  
390 flexibility to model complex, non-linear dependencies but lack the parametric guarantees of  
391 convergence.

392 To explore the predictive power of combining machine learning classifiers as feature embed-  
393 ding tools, we use the Extreme Gradient Boosting Survival Embedding (XGBSE) model<sup>16</sup>,  
394 which combines an Accelerated Failure Time framework as tree-based embedding with logistic  
395 regressions to estimate bin-level survival probabilities<sup>36</sup>.

396 As a control, we also evaluate predictions using classification models, including LightGBM<sup>26</sup>,  
397 LightGBM with monotone risk increase for activity duration, logistic regression with Light-  
398 GBM tree embeddings, k-Nearest Neighbors<sup>27;37</sup>, and logistic regression<sup>28;37</sup>. All classification  
399 models incorporate Bayesian hyperparameter optimization<sup>38</sup>.

## 400 Calculating Injury Probabilities Based on Survival Curves

401 All survival-based approaches estimate a survival function  $S_\theta : (0, \infty] \rightarrow [0, 1]$ , where the  
402 parameterization  $\theta$  depends on the specific method employed. In Cox-based models, such as  
403 Coxnet<sup>18</sup>,  $\theta \in \mathbb{R}^{n_\theta}$  represents the coefficients of the proportional hazards model. Tree-based  
404 methods, including Accelerated Failure Time<sup>17</sup> and Boosted-Cox<sup>30</sup>, parameterize  $\theta \in \mathbb{F}$ ,  
405 where  $\mathbb{F}$  denotes the space of additive tree ensembles. The XGBSE model<sup>16</sup> combines these  
406 approaches, with  $\theta \in \mathbb{F} \times \mathbb{R}^{n_\theta}$ , comprising a tree ensemble for the underlying Accelerated Fail-  
407 ure Time model and real-valued parameters for logistic regressions used in post-processing.

408 Given a parameterization  $\theta$ , the survival function  $S_\theta(m|x)$ , conditioned on covariates  $x \in \mathbb{R}^{n_x}$ ,  
409 takes as input the number of player minutes  $m$  and outputs the probability that an individual  
410 does not sustain an injury up to minute  $m$ . The probability of sustaining an injury up to  
411 minute  $m$  is expressed as  $\mathbb{P}_\theta(m|x) = 1 - S_\theta(m|x)$ .

412 The monotonicity of the survival function ensures that the probability of injury increases as  
413 the number of minutes played increases. As our data reports only the actual minutes played,  
414 which serve as lower bounds for planned activity durations on injury days, we compute  $p_{i,t}$  on  
415 injury days as the mean injury probability for all training durations of players who trained  
416 at least as long as the injured player on that day. For non-injury observations, we use the  
417 recorded time of activity.

418 Given that readers may not be familiar with the XGBSE model, we provide a detailed  
419 explanation of its approach to estimating survival curves in the supplementary material.

## 420 Probability Calibration with Beta Calibration

421 To adjust the estimated probabilities to align with observed injury frequencies, we employed  
422 probability calibration for all models. Using the model parameters  $\hat{\theta}$  estimated from the  
423 training set  $\mathcal{D}^{\text{Train}}$ , we calculated the uncalibrated probabilities  $\hat{p}_{i,t} = \mathbb{P}_{\hat{\theta}}(m|x)$  for each ob-  
424 servation in the calibration set  $\mathcal{D}^{\text{Cal}}$ , consisting of unseen data points  $(y_{i,t}, x_{i,t}, m_{i,t})$ .

425 The uncalibrated probabilities  $\hat{p}_{i,t}$  were adjusted using beta calibration<sup>19</sup>, which applies the  
426 following mapping function:

$$\hat{p}_{i,t}^c = f_{\alpha}(\hat{p}_{i,t}) = \left( 1 + \frac{(1 - \hat{p}_{i,t})^{\alpha_2}}{\exp(\alpha_3) \cdot \hat{p}_{i,t}^{\alpha_1}} \right)^{-1}, \quad (1)$$

427 where  $\alpha = (\alpha_1, \alpha_2, \alpha_3)'$  are the parameters of the calibration function. These parameters  
428 were estimated by fitting the calibration function to the calibration set  $\mathcal{D}^{\text{Cal}}$ .

## 429 Day Importance Weighting

430 To account for the varying importance of a player missing a particular day, we introduced  
431 the concept of match importance. Match importance assigns an individual cost  $c_{i,t} \in \mathbb{R}_+$  to  
432 each day  $t$  for a player  $i$ , reflecting the impact of their absence.

433 These costs are specified by decision-makers and can be dynamically adjusted to reflect  
434 context-specific priorities, such as the higher importance of cup knockout matches compared  
435 to practice sessions. As the definition of these costs is highly club-specific, we distinguished  
436 for our analysis between the cost of missing a match  $c^{(G)} \in \mathbb{R}_+$  and the cost of missing a  
437 practice  $c^{(P)} \in \mathbb{R}_+$ , both of which are assumed to be constant across individuals and days.

438 Since only the relative values of  $c^{(G)}$  and  $c^{(P)}$  are relevant for our analysis, we defined the  
439 relative match importance ratio  $R = c^{(G)}/c^{(P)}$ . This ratio is used to transform practice days  
440 into match-weighted practice days by dividing the number of practice days by  $R$ . In an  
441 application of a club the ratio  $R$  can be day and individual specific, allowing individualized  
442 tailored weighting of days for players and competitions.

## 443 Cost of Playing and Resting

444 The decision for a player to participate in a match or practice session carries a potential  
445 cost, depending on whether the player sustains an injury. To quantify this, we defined a  
446 cost function incorporating the predicted injury label  $\hat{y}_{i,t} \in \{0, 1\}$ , the actual injury outcome  
447  $y_{i,t} \in \{0, 1\}$ , the number of missed matches  $D_{i,t}^G \in \mathbb{N}$  and practices  $D_{i,t}^P \in \mathbb{N}$ , the match  
448 indicator  $G_{i,t}$  (1 for match days, 0 otherwise), and the match importance ratio  $R$ .

449 If a player rests, the cost reflects the missed day, with a value of 1 for match days and  $1/R$   
 450 for practice days. If a player participates and sustains an injury, the cost accounts for the  
 451 missed matches and practices, expressed in match-weighted days. To consolidate these cases,  
 452 we defined the overall cost function

$$C(y_{i,t}, \hat{y}_{i,t}) = \begin{cases} \mathbb{1}_{\{G_{i,t}=1\}} + \frac{\mathbb{1}_{\{G_{i,t}=0\}}}{R} & \text{if } \hat{y}_{i,t} = 1, \\ D_{i,t}^G + \frac{D_{i,t}^P}{R} & \text{if } \hat{y}_{i,t} = 0, y_{i,t} = 1, \\ 0 & \text{if } \hat{y}_{i,t} = 0, y_{i,t} = 0. \end{cases} \quad (2)$$

453 This cost function provides a unified framework for quantifying decision outcomes and facil-  
 454 itates the calculation of conditionally expected costs, discussed in the next subsection.

## 455 Conditionally Expected Costs Under Injury Duration Uncertainty

456 The cost function  $C(y_{i,t}, \hat{y}_{i,t})$  depends on the actual injury outcome  $y_{i,t}$ , the number of missed  
 457 matches  $D_{i,t}^G$ , and missed practices  $D_{i,t}^P$  in the case of an injury. However, a priori, the injury  
 458 outcome and injury durations are unknown.

459 Using the estimated calibrated probabilities, the expectation of the costs, conditioned on  $D_{i,t}^G$   
 460 and  $D_{i,t}^P$ , is dependent on the decision to rest a player ( $\hat{y}_{i,t}$ ) and is expressed as:

$$\mathbb{E}(C(y_{i,t}, \hat{y}_{i,t}) | D_{i,t}^G, D_{i,t}^P) = \begin{cases} \mathbb{1}_{\{G_{i,t}=1\}} + \frac{\mathbb{1}_{\{G_{i,t}=0\}}}{R} & \text{if } \hat{y}_{i,t} = 1, \\ \hat{p}_{i,t}^c \cdot \left[ D_{i,t}^G + \frac{D_{i,t}^P}{R} \right] & \text{if } \hat{y}_{i,t} = 0. \end{cases} \quad (3)$$

461 If  $\hat{y}_{i,t} = 1$ , the conditionally expected costs incorporate stochasticity due to the uncertainty  
 462 in the injury outcome. The player incurs a cost associated with an injury, weighted by the  
 463 calibrated injury probability  $\hat{p}_{i,t}^c$ , or no cost if no injury occurs with probability  $1 - \hat{p}_{i,t}^c$ . If  
 464  $\hat{y}_{i,t} = 0$ , the costs are deterministic because the player rests and misses the session, and the  
 465 conditionally expected costs are equal to the cost function for resting.

466 Since  $D_{i,t}^G$  and  $D_{i,t}^P$  are unknown a priori, the conditionally expected cost  $\mathbb{E}(C(y_{i,t}, \hat{y}_{i,t}) | D_{i,t}^G, D_{i,t}^P)$   
 467 is itself a random variable that depends on the underlying distribution of  $D_{i,t}^G$  and  $D_{i,t}^P$ .

## 468 Augmenting Data Set with Historical Injury Durations

469 The distribution of missed matches and practices due to injury is critical for deriving the  
 470 distribution of conditionally expected costs. To estimate this, we utilized historical data and  
 471 augmented it with additional data from SoccerDonna<sup>39</sup>.

472 Our dataset initially contained 83 recorded non-contact injuries. To improve the robustness  
 473 of the injury duration distribution estimate, we included injury data from players listed on

474 the first teams of the top female football divisions in Spain (Primera División de la Liga  
 475 de Fútbol Femenino), Germany (Frauen-Bundesliga), England (Women’s Super League),  
 476 and the United States (National Women’s Soccer League) during the 2022/2023 season.  
 477 This augmented dataset comprised injury durations  $D_j$  for  $j = 1, 2, \dots, 494$ , significantly  
 478 increasing the sample size by 411 observations. We included all injuries indicative of non-  
 479 contact mechanisms; a complete list of the selection criteria is provided in the Supplementary  
 480 Table S1.

481 To estimate the number of missed matches and practices from total injury durations, we  
 482 applied a linear regression approach to partition the injury durations into match days and  
 483 practice days. Specifically, we fitted the following regression, once for both  $k = G, P$ , on the  
 484 combined training and calibration:

$$D_{i,t}^k = \beta_0^k + \beta_1^k D_{i,t} + \varepsilon_{i,t}^k, \quad (4)$$

485 where  $D_{i,t}$  is the total injury duration recorded for individual  $i$  on day  $t$ . The fitted parameters  
 486  $\hat{\beta}_0^G, \hat{\beta}_1^G, \hat{\beta}_0^P, \hat{\beta}_1^P$  were then applied to the augmented injury dataset to estimate the number  
 487 of missed matches ( $D_j^G$ ) and practices ( $D_j^P$ ) for each injury observation. The distribution of  
 488  $D_j$ , as well as the fitted linear regression lines for the entire dataset, are provided within the  
 489 Supplementary Figure S8.

## 490 Decision Thresholds and Rest Benefit Certainty

491 Given the calibrated probabilities  $\hat{p}_{i,t}^c$ , the cost function  $C(y_{i,t}, \hat{y}_{i,t})$ , and the joint distribution  
 492 of missed matches  $D_{i,t}^G$  and practices  $D_{i,t}^P$ , decision-makers must determine whether a player  
 493  $i$  should play or rest on day  $t$ . This decision is based on whether the expected cost of playing  
 494 exceeds the expected cost of resting:

$$\text{Rest Player if: } \underbrace{\mathbb{E}(C(y_{i,t}, \hat{y}_{i,t} = 1) | D_{i,t}^G, D_{i,t}^P)}_{\text{Expected Cost of Playing}} > \underbrace{\mathbb{E}(C(y_{i,t}, \hat{y}_{i,t} = 0) | D_{i,t}^G, D_{i,t}^P)}_{\text{Expected Cost of Resting}}. \quad (5)$$

495 Using the previously derived formula for expected costs, this inequality can be reformulated  
 496 in terms of the calibrated probabilities  $\hat{p}_{i,t}^c$ , the match indicator  $\mathbb{1}_{\{G_{i,t}=1\}}$ , and the match  
 497 valuation parameter  $R$ . Solving for  $\hat{p}_{i,t}^c$  provides the following decision thresholds

$$\begin{aligned} \text{Match Day: } \hat{p}_{i,t}^c &> \left[ D_{i,t}^G + \frac{D_{i,t}^P}{R} \right]^{-1}, \\ \text{Practice Day: } \hat{p}_{i,t}^c &> \frac{1}{R} \left[ D_{i,t}^G + \frac{D_{i,t}^P}{R} \right]^{-1}. \end{aligned} \quad (6)$$

498 For match days, it is optimal to rest a player if the probability of injury exceeds the inverse  
 499 of the match-weighted injury duration. For practice days, this threshold is further adjusted  
 500 by the inverse match valuation parameter  $R$ . This adjustment reflects that, all else being  
 501 equal, the probability of injury must be  $R$  times higher on match days than on practice days  
 502 to justify resting the player.

503 Using the joint distribution of  $D_{i,t}^G$  and  $D_{i,t}^P$ , the probability  $Q_{i,t}$  that the threshold condition  
 504 in Inequality 6 is satisfied can be computed. For match days, this probability is given by:

$$Q_{i,t} = \mathbb{P} \left( \hat{p}_{i,t}^c > \left[ D_{i,t}^G + \frac{D_{i,t}^P}{R} \right]^{-1} \right). \quad (7)$$

505 The value  $Q_{i,t}$  represents the probability that resting a player is the optimal decision. Decision-  
 506 makers specify a threshold level of certainty, referred to as the Rest Benefit Certainty (RBC).  
 507 If  $Q_{i,t} > \text{RBC}$ , the player is rested; otherwise, she plays.

## 508 Availability Gain Evaluation Metric

509 The Availability Gain (AG) metric evaluates decision-making strategies by balancing the  
 510 negative impact of unnecessary rests with the positive impact of detecting injuries that  
 511 result in missed match and practice days. Days are weighted as match days using the match  
 512 valuation parameter  $R$ .

513 The metric assigns a negative value to rested days, regardless of whether an injury occurs,  
 514 and a positive value to detected injuries. It is defined as:

$$AG(y_{i,t}, \hat{y}_{i,t}) = \begin{cases} -\mathbb{1}_{\{G_{i,t}=1\}} - \frac{\mathbb{1}_{\{G_{i,t}=0\}}}{R} & \text{if } \hat{y}_{i,t} = 1, y_{i,t} = 0, \\ -\mathbb{1}_{\{G_{i,t}=1\}} - \frac{\mathbb{1}_{\{G_{i,t}=0\}}}{R} + D_{i,t}^G + \frac{D_{i,t}^P}{R} & \text{if } \hat{y}_{i,t} = 1, y_{i,t} = 1, \\ 0 & \text{if } \hat{y}_{i,t} = 0. \end{cases} \quad (8)$$

515 In this formulation, the negative term corresponds to the cost of resting, and the positive  
 516 term accounts for the benefit of detecting injuries in terms of missed match ( $D_{i,t}^G$ ) and practice  
 517 days ( $D_{i,t}^P$ ). The AG metric over the test set is computed as the sum of availability gains  
 518 across all test set observations.

## 519 Interpretation of Mean Average Precision and Area Under the Re- 520 ceiving Operator Curve

521 We quantified discrimination using the Mean Average Precision (MAP) and the Area Under  
 522 the Receiver Operating Characteristic Curve (AUROC). MAP evaluates how well a model

523 balances precision and recall, making it particularly suited for imbalanced datasets because it  
524 penalizes false positives and rewards correct identification of the minority class. In contrast,  
525 AUROC represents the average sensitivity and specificity across varying decision thresholds,  
526 providing a broader assessment of performance.

## 527 **Injury Data Recording**

528 The medical team managed and documented injuries using a validated electronic medical  
529 record system (COR version 2.0; FCB). Injury assessments were performed by the team's  
530 medical physician, in conjunction with the FC Barcelona medical department, following  
531 standardized diagnosis and return-to-play protocols as outlined in the club's guidelines<sup>24</sup>. The  
532 analysis focused on non-contact injuries affecting muscles, tendons, ligaments, and cartilage.

## 533 **Recording of Workload Features**

534 The external workload during training sessions and matches was monitored using GPS data  
535 obtained from the WIMU PRO™ device (RealtrackSystems S.L., Almeria, Spain)<sup>40</sup>. The col-  
536 lected data underwent pre-processing through SPRO™ Software (version 927) (RealtrackSys-  
537 tems, Almeria, Spain). This software compiled the data in RAW format and facilitated the  
538 generation of continuous training-load variables. The GPS training-load variables used in  
539 this study for include the distance covered in meters and the high metabolic load distance  
540 (HMLD; measured in meters, defined as the distance covered when metabolic power exceeds  
541 25.5 Watt/kg, at speeds above 5.5 m/s). Additionally, the number of accelerations (higher  
542 than 3 m/s<sup>2</sup>) and decelerations, along with the respective distances in meters covered dur-  
543 ing accelerations as well as decelerations, are considered. Finally, the distance covered while  
544 high-speed running per minute (relative HSR) is included, defined as running at speeds above  
545 18 km/h.

## 546 **Feature Engineering**

547 Baseline features included per-minute training load variables derived from player tracking  
548 data. For the classification comparison models, the duration in minutes was added as an  
549 additional feature. To capture training load dynamics, new variables were designed following  
550 the methodology of Rossi et al.<sup>6</sup>. These included exponentially weighted moving averages,  
551 with 21 features calculated to represent past training loads. The acute chronic workload ratio  
552 was defined as the ratio of the load over the last six days to the load over the last 28 days,  
553 assessing shifts between acute and chronic training loads. The monotonic workload ratio  
554 was calculated as the mean load over the last seven days divided by the standard deviation

555 during the same period, capturing variability in recent training loads. Additional features  
556 incorporated into the models were age, treated as a continuous variable, and a binary match  
557 indicator variable.

558 The list of all used features and their computations, if applicable, are given within the  
559 supplementary Table S3.

## 560 **Feature Pre-Processing**

561 For numerical features used in the Coxnet, k-Nearest Neighbors, and Logistic Regression mod-  
562 els, standardization was performed using the mean and standard deviation of each feature.  
563 Principal Component Analysis (PCA) was then applied to reduce dimensionality, retaining  
564 components that captured 95% of the total variance in the data.

## 565 **Hyperparameter Optimization for Classification Models**

566 Hyperparameter optimization for the standard machine learning models was carried out  
567 using Bayesian hyperparameter optimization<sup>38;41</sup> with 60 iterations, which is standard in  
568 the literature<sup>42</sup>. The final parameters and details on the hyperparameter search space are  
569 provided in the Supplementary Material.

## 570 **Implementation and Code Availability**

571 All models were implemented in Python 3.11.8. The survival-based models were implemented  
572 as follows: the XGBSE model utilized the xgbse package version 0.3.3, the Accelerated Fail-  
573 ure Time model was implemented with the xgboost package version 2.1.1, and the Coxnet  
574 and CoxBoost models were developed using the scikit-survival package version 0.23.0. The  
575 classification comparison models, including LightGBM, LightGBM with monotone duration  
576 of activity, k-Nearest Neighbors, and Logistic Regression, were developed using scikit-learn  
577 version 1.5.2. The LightGBM Tree Embedding model was implemented using the same Light-  
578 GBM and Logistic Regression setup. Code will be made available from the corresponding  
579 author upon reasonable request.

## 580 **Data Availability**

581 Data will be available upon reasonable request, considering ethical and privacy concerns,  
582 from the author Gil Rodas.

## 583 **Ethical Considerations**

584 This study adhered to the Declaration of Helsinki guidelines<sup>43</sup> and received approval from  
585 both the Barça Innovation Hub's local committee and the Ethics Committee of Consell  
586 Català de l'Esport (code 012/CEICGC/2021). Participants were informed about the study's  
587 potential risks and benefits. All personal data and results were anonymized to maintain con-  
588 fidentiality and comply with the criteria specified by the General Data Protection Regulation  
589 (GDPR).

## 590 **Funding**

591 This study was funded by the Ministry of Research and Universities, Government of Catalo-  
592 nia, under reference AGAUR 2023 PROD 00020. It also received financial support by the  
593 German Research Foundation (Deutsche Forschungsgemeinschaft, DFG) under Germany's  
594 Excellence Strategy (EXC 2047 - 390685813 and EXC 2151 - 390873048), the University of  
595 Bonn (via the Schlegel Professorship of JH), and by the European Union. Views and opinions  
596 expressed are however those of the author(s) only and do not necessarily reflect those of the  
597 European Union or the European Research Council Executive Agency. Neither the European  
598 Union nor the granting authority can be held responsible for them. This work is supported  
599 by ERC grant INTEGRATE, grant agreement number 101126146.

600 The funders had no role in the study design, data collection, data analyses, data interpreta-  
601 tion, writing, or submission of this manuscript.

## 602 **Author Contributions**

603 M.H. developed the statistical and mathematical theory and implemented the models in  
604 Python. B.C. conducted robustness checks. E.F., X. Y. and G.R. collected the data. M.H.  
605 visualized the results. M.H., J.H., and J.R.G. conceptualized the study. M.H., J.H., and  
606 J.R.G. wrote the manuscript. J.H. and J.R.G. oversaw the research process. All authors  
607 reviewed and approved the final version.

## 608 **Declaration of Interests**

609 E.F, and G.R. are employed by Barça Innovation Hub. J.R.G. serves as a scientific advisor to  
610 Made of Genes, for which he receives financial compensation. B.C. is pursuing an industrial

611 PhD in collaboration with ISGlobal and Made of Genes. All other authors have no competing  
612 interests to declare.

## 613 References

- 614 [1] Martin Hägglund, Markus Waldén, Henrik Magnusson, Karolina Kristenson, Håkan  
615 Bengtsson, and Jan Ekstrand. Injuries affect team performance negatively in profes-  
616 sional football: an 11-year follow-up of the ufa champions league injury study. *British*  
617 *journal of sports medicine*, 47(12):738–742, 2013.
- 618 [2] Eyal Eliakim, Elia Morgulev, Ronnie Lidor, and Yoav Meckel. Estimation of injury costs:  
619 financial damage of english premier league teams’ underachievement due to injuries. *BMJ*  
620 *Open Sport & Exercise Medicine*, 6(1):e000675, 2020.
- 621 [3] Olivia A Hurley. Impact of player injuries on teams’ mental states, and subsequent  
622 performances, at the rugby world cup 2015. *Frontiers in psychology*, 7:807, 2016.
- 623 [4] SR Filbay, AG Culvenor, IN Ackerman, TG Russell, and KM Crossley. Quality of life in  
624 anterior cruciate ligament-deficient individuals: a systematic review and meta-analysis.  
625 *British Journal of Sports Medicine*, 49(16):1033–1041, 2015.
- 626 [5] Howden. Howden injury index: Sports injury trends and insights, 2023. Available at:  
627 <https://www.howdengroup.com>.
- 628 [6] Alessio Rossi, Luca Pappalardo, Paolo Cintia, F Marcello Iaia, Javier Fernández, and  
629 Daniel Medina. Effective injury forecasting in soccer with gps training data and machine  
630 learning. *PLOS One*, 13(7):e0201264, 2018.
- 631 [7] Alejandro López-Valenciano, Francisco Ayala, José Miguel Puerta, Mark De Ste Croix,  
632 Francisco Vera-García, Sergio Hernández-Sánchez, Iñaki Ruiz-Pérez, and Gregory Myer.  
633 A preventive model for muscle injuries: a novel approach based on learning algorithms.  
634 *Medicine and science in sports and exercise*, 50(5):915, 2018.
- 635 [8] Francisco Ayala, Alejandro López-Valenciano, Jose Antonio Gámez Martín, Mark De Ste  
636 Croix, Francisco J Vera-Garcia, Maria del Pilar Garcia-Vaquero, Iñaki Ruiz-Pérez, and  
637 Gregory D Myer. A preventive model for hamstring injuries in professional soccer:  
638 Learning algorithms. *International journal of sports medicine*, 40(05):344–353, 2019.
- 639 [9] Nikki Rommers, Roland Rössler, Evert Verhagen, Florian Vandecasteele, Steven Ver-  
640 stockt, Roel Vaeyens, Matthieu Lenoir, Eva D’Hondt, and Erik Witvrouw. A machine  
641 learning approach to assess injury risk in elite youth football players. *Medicine and*  
642 *science in sports and exercise*, 52(8):1745–1751, 2020.

- 643 [10] Garrett S Bullock, Joseph Mylott, Tom Hughes, Kristen F Nicholson, Richard D Riley,  
644 and Gary S Collins. Just how confident can we be in predicting sports injuries? a system-  
645 atic review of the methodological conduct and performance of existing musculoskeletal  
646 injury prediction models in sport. *Sports medicine*, 52(10):2469–2482, 2022.
- 647 [11] Christopher Leckey, Nicol van Dyk, Cailbhe Doherty, Aonghus Lawlor, and Eamonn  
648 Delahunt. Machine learning approaches to injury risk prediction in sport: a scoping  
649 review with evidence synthesis. *British Journal of Sports Medicine*, 2024.
- 650 [12] Jon L Oliver, Francisco Ayala, Mark BA De Ste Croix, Rhodri S Lloyd, Greg D Myer,  
651 and Paul J Read. Using machine learning to improve our understanding of injury risk  
652 and prediction in elite male youth football players. *Journal of science and medicine in  
653 sport*, 23(11):1044–1048, 2020.
- 654 [13] Emmanuel Vallance, Nicolas Sutton-Charani, Abdelhak Imoussaten, Jacky Montmain,  
655 and Stéphane Perrey. Combining internal-and external-training-loads to predict non-  
656 contact injuries in soccer. *Applied Sciences*, 10(15):5261, 2020.
- 657 [14] Ewout W Steyerberg, Andrew J Vickers, Nancy R Cook, Thomas Gerds, Mithat Go-  
658 nen, Nancy Obuchowski, Michael J Pencina, and Michael W Kattan. Assessing the  
659 performance of prediction models: a framework for traditional and novel measures. *Epi-  
660 demiology*, 21(1):128–138, 2010.
- 661 [15] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern  
662 neural networks. In *International conference on machine learning*, pages 1321–1330.  
663 PMLR, 2017.
- 664 [16] Davi Vieira, Gabriel Gimenez, Guilherme Marmerola, and Vitor Estima. Xgboost sur-  
665 vival embeddings: improving statistical properties of xgboost survival analysis imple-  
666 mentation, 2021. URL <https://loft-br.github.io/xgboost-survival-embeddings/>.
- 667 [17] Avinash Barnwal, Hyunsu Cho, and Toby Hocking. Survival regression with accelerated  
668 failure time model in xgboost. *Journal of Computational and Graphical Statistics*, 31  
669 (4):1292–1302, 2022.
- 670 [18] Noah Simon, Jerome H Friedman, Trevor Hastie, and Rob Tibshirani. Regularization  
671 paths for cox’s proportional hazards model via coordinate descent. *Journal of statistical  
672 software*, 39:1–13, 2011.
- 673 [19] Meelis Kull, Telmo Silva Filho, and Peter Flach. Beta calibration: a well-founded and  
674 easily implemented improvement on logistic calibration for binary classifiers. In *Artificial  
675 intelligence and statistics*, pages 623–631. PMLR, 2017.

- 676 [20] Monika E Maros, David Capper, David T W Jones, Volker Hovestadt, Andreas von  
677 Deimling, Stefan M Pfister, Martin Sill, and Felix Sahm. Machine-learning workflows  
678 for the classification of dna methylation array data. *Nature Protocols*, 15(2):479–512,  
679 2019. doi: 10.1038/s41596-019-0251-6.
- 680 [21] Haoyue Huang, Jiawei Zheng, Hao Zhang, Bing Wang, Mingxi Zhou, Jiacheng Zhang,  
681 and Ying Wang. Machine learning-based tissue of origin classification for cancers of  
682 unknown primary. *Nature Communications*, 13(1):4327, 2022. doi: 10.1038/s41467-022-  
683 31666-w.
- 684 [22] Gil Rodas, Lourdes Osaba, David Arteta, Ricard Pruna, Dolors Fernández, and Alejan-  
685 dro Lucia. Genomic prediction of tendinopathy risk in elite team sports. *International*  
686 *Journal of Sports Physiology and Performance*, 15(4):489–495, 2019.
- 687 [23] Juan Ramon González, Alejandro Cáceres, Eva Ferrer, Laura Balagué-Dobón, Xavier  
688 Escribà-Montagut, David Sarrat, Guillermo Quintás, and Rodas Rodas. Predicting  
689 injuries in elite female football playe... *International Journal of Sports Physiology and*  
690 *Performance*, 19:661–669, 2024.
- 691 [24] Ricard Pruna, Thor Einar Andersen, Ben Clarsen, and Alan McCall. Muscle injury  
692 guide: Prevention of and return to play from muscle injuries. Technical report, Barca  
693 Innovation Hub, 2018.
- 694 [25] Jill Cook, Gil Rodas, Alan McCall, Ricard Pruna, Rochelle Kennedy, and Lluís Til.  
695 *Tendon Injuries in Football Players: FC Barcelona 2021 Tendon Guide*. FC Barcelona:  
696 Barça Innovation Hub, Barcelona, Spain, 1st edition, 2021.
- 697 [26] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye,  
698 and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. *Advances*  
699 *in Neural Information Processing Systems*, 30, 2017.
- 700 [27] Jorma Laaksonen and Erkki Oja. Classification with learning k-nearest neighbors. In  
701 *Proceedings of international conference on neural networks (ICNN'96)*, volume 3, pages  
702 1480–1483. IEEE, 1996.
- 703 [28] David W Hosmer Jr, Stanley Lemeshow, and Rodney X Sturdivant. *Applied logistic*  
704 *regression*. John Wiley & Sons, 2013.
- 705 [29] Ahmed Naglah, Fahmi Khalifa, Ali Mahmoud, Mohammad Ghazal, Paul Jones, Teena  
706 Murray, Adel S Elmaghraby, and Ayman El-Baz. Athlete-customized injury prediction  
707 using training load statistical records and machine learning. In *2018 IEEE international*  
708 *symposium on signal processing and information technology (ISSPIT)*, pages 459–464.  
709 IEEE, 2018.

- 710 [30] Greg Ridgeway. The state of boosting. *Computing science and statistics*, pages 172–181,  
711 1999.
- 712 [31] Toshihito Takahashi, Kazunori Nozaki, Tomoya Gonda, Tomoaki Mameno, and  
713 Kazunori Ikebe. Deep learning-based detection of dental prostheses and restorations.  
714 *Scientific Reports*, 11(1):1960, 2021.
- 715 [32] Todd Hollon, Cheng Jiang, Asadur Chowdury, Mustafa Nasir-Moin, Akhil Kondepudi,  
716 Alexander Aabedi, Arjun Adapa, Wajd Al-Holou, Jason Heth, Oren Sagher, et al.  
717 Artificial-intelligence-based molecular classification of diffuse gliomas using rapid, label-  
718 free optical imaging. *Nature medicine*, 29(4):828–832, 2023.
- 719 [33] Cox R David et al. Regression models and life tables (with discussion). *Journal of the*  
720 *Royal Statistical Society*, 34(2):187–220, 1972.
- 721 [34] Sebastian Pölsterl. scikit-survival: A library for time-to-event analysis built on top of  
722 scikit-learn. *Journal of Machine Learning Research*, 21(212):1–6, 2020.
- 723 [35] Tianqi Chen and Carlos Guestrin. XGBoost: A scalable tree boosting system. In *Pro-*  
724 *ceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery*  
725 *and Data Mining*, KDD '16, pages 785–794, New York, NY, USA, 2016. ACM. ISBN  
726 978-1-4503-4232-2.
- 727 [36] Xinran He, Junfeng Pan, Ou Jin, Tianbing Xu, Bo Liu, Tao Xu, Yanxin Shi, Antoine  
728 Atallah, Ralf Herbrich, Stuart Bowers, et al. Practical lessons from predicting clicks on  
729 ads at facebook. In *Proceedings of the eighth international workshop on data mining for*  
730 *online advertising*, pages 1–9, 2014.
- 731 [37] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blon-  
732 del, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau,  
733 M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python.  
734 *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- 735 [38] Tim Head, Manoj Kumar, Holger Nahrstaedt, Gilles Louppe, and Iaroslav Shcherbatyi.  
736 scikit-optimize/scikit-optimize, October 2021.
- 737 [39] Soccerdonna. Soccerdonna - women’s soccer database. <https://www.soccerdonna.de>.
- 738 [40] Marc Guitart, Martí Casals, David Casamichana, Jordi Cortés, Francesc Xavier Valle,  
739 Alan McCall, Francesc Cos, and Gil Rodas. Use of gps to measure external load and  
740 estimate the incidence of muscle injuries in men’s football: A novel descriptive study.  
741 *PLOS One*, 17(2):e0263494, 2022.

- 742 [41] Donald R. Jones, Matthias Schonlau, and William J. Welch. Efficient global optimization  
743 of expensive black-box functions. *Journal of Global Optimization*, 13(4):455–492, 1998.  
744 doi: 10.1023/A:1008306431147.
- 745 [42] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization  
746 of machine learning algorithms. *Advances in neural information processing systems*, 25,  
747 2012.
- 748 [43] World Medical Association. Wma declaration of helsinki – ethical princi-  
749 ples for medical research involving human subjects, 2013. Retrieved from  
750 [https://www.wma.net/policies-post/wma-declaration-of-helsinki-ethical-](https://www.wma.net/policies-post/wma-declaration-of-helsinki-ethical-principles-for-medical-research-involving-human-subjects/)  
751 [principles-for-medical-research-involving-human-subjects/](https://www.wma.net/policies-post/wma-declaration-of-helsinki-ethical-principles-for-medical-research-involving-human-subjects/).

# 752 **Supplementary Material to: Informed Injury Prediction** 753 **in Elite Football: Decision Theory meets Machine** 754 **Learning**

755 Manuel Huth<sup>1,2</sup>, Berta Canal-Simón<sup>3,4</sup>, Eva Ferrer<sup>5,6</sup>, Gil Rodas<sup>5,6,7</sup>,  
756 Xavier Yanguas<sup>5</sup>, Jan Hasenauer<sup>1,2†</sup>, and Juan R González<sup>3,8,9†\*</sup>

757 <sup>1</sup> Life and Medical Sciences (LIMES) Institute, University of Bonn, Bonn, Germany

758 <sup>2</sup> Bonn Center for Mathematical Life Sciences, University of Bonn, Bonn, Germany

759 <sup>3</sup> Barcelona Institute for Global Health (ISGlobal), Barcelona, Spain.

760 <sup>4</sup> Made of Genes, Barcelona, Spain.

761 <sup>5</sup> Medical Department of Football Club Barcelona (FIFA Medical Centre of Excellence), and Barça  
762 Innovation Hub of Football Club Barcelona, Barcelona, Spain.

763 <sup>6</sup> Sports and Exercise Medicine Unit, Hospital Clinic and Sant Joan de Déu, Barcelona, Spain.

764 <sup>7</sup> Leitat technological Center, Terrassa, Spain.

765 <sup>8</sup> CIBER in Epidemiology and Public Health (CIBERESP), Barcelona, Spain.

766 <sup>9</sup> Department of Mathematics, Universitat Autònoma de Barcelona (UAB), Barcelona, Spain.

767 † Authors contributed equally; \* Correspondence: [juan.gonzalez@isglobal.org](mailto:juan.gonzalez@isglobal.org)

768 This supplementary material provides additional technical details, illustrative examples, and method-  
769 ological insights that complement our main study, Informed Injury Prediction in Elite Football:  
770 Decision Theory meets Machine Learning. We begin by presenting detailed examples of individual  
771 decision-making, where we contrast the expected costs of resting versus playing under various injury  
772 scenarios and illustrate the concept of Rest Benefit Certainty (RBC). Next, we showcase the avail-  
773 ability gain metric to quantify how different injury durations affect player availability, incorporating  
774 the impact of match valuation.

775 Further, we describe the Extreme Gradient Boosting Survival Embedding (XGBSE) model used  
776 for estimating survival curves from censored injury data, outlining its three-step procedure—from  
777 fitting a censored Accelerated Failure Time model to feature space transformation and logistic  
778 regression on binned survival times. Finally, we explain our imputation strategies for integrating  
779 national team training exposure into the overall training history, ensuring that gaps in the data are  
780 appropriately addressed.

## 781 **Illustrative Example: Individual Decision-Making**

782 Consider a scenario where decision-makers must decide whether a player should rest or participate  
783 in training the next day. The model predicts a 5% probability of injury for the upcoming day, which  
784 is a practice day. Injuries are assumed to last, with equal probability, 7, 14, or 21 days. Matches  
785 occur once per week, and matches are valued  $R = 5$  times more than practices.

### 786 **Expected Cost of Resting vs. Playing**

787 The cost of resting the player on a practice day is given by the match-weighted value of a practice  
788 day:

$$\text{Cost of Resting} = \frac{1}{R} = \frac{1}{5} = 0.2. \quad (9)$$

789 The expected cost of playing, conditioned on the injury duration, is calculated as follows:

$$\text{Injury Probability} \cdot (\text{Number of Practice Days}/R + \text{Number of Match Days}), \quad (10)$$

790 such that

$$\begin{aligned} \text{(a) Injury lasts 7 days: } & 0.05 \cdot \left( \frac{6}{5} + 1 \right) = 0.11, \\ \text{(b) Injury lasts 14 days: } & 0.05 \cdot \left( \frac{12}{5} + 2 \right) = 0.22, \\ \text{(c) Injury lasts 21 days: } & 0.05 \cdot \left( \frac{18}{5} + 3 \right) = 0.33. \end{aligned} \quad (11)$$

791 Comparing the resting cost of 0.2 to the conditional expected costs of playing, we observe that  
792 only in scenario (a), where the injury lasts 7 days, playing yields a lower cost. In 66.67% of cases  
793 (scenarios b and c), playing has a higher expected cost, making resting the favorable decision under  
794 these assumptions.

### 795 **Rest Benefit Certainty (RBC)**

796 Decision-makers may require a minimum level of certainty, referred to as Rest Benefit Certainty  
797 (RBC), to rest the player. For instance, if the required certainty level is 50%, resting is optimal  
798 because playing has a higher expected cost in 66.67% of cases. However, if the required certainty is  
799 70%, decision-makers would not rest the player. The RBC represents the required confidence level  
800 to justify resting but does not directly affect the expected costs.

### 801 **Influence of Match Valuation**

802 The match valuation parameter  $R$  directly influences the costs. If matches are valued only 2 times  
803 more than practices ( $R = 2$ ), the cost of resting becomes:

$$\text{Cost of Resting} = \frac{1}{R} = \frac{1}{2} = 0.5. \quad (12)$$

804 The expected costs of playing are recalculated as:

$$\begin{aligned} \text{(a) Injury lasts 7 days: } & 0.05 \cdot \left( \frac{6}{2} + 1 \right) = 0.2, \\ \text{(b) Injury lasts 14 days: } & 0.05 \cdot \left( \frac{12}{2} + 2 \right) = 0.4, \\ \text{(c) Injury lasts 21 days: } & 0.05 \cdot \left( \frac{18}{2} + 3 \right) = 0.6. \end{aligned} \quad (13)$$

805 Comparing these values to the new resting cost of 0.5, we see that only in scenario (c), where  
806 the injury lasts 21 days, resting yields a lower cost. As this outcome has a probability of 33.33%,  
807 decision-makers with a RBC of 50% would now decide to let the player participate.

808 This example highlights how match valuation influences decision-making. When matches are highly  
809 valued ( $R = 5$ ), even short injuries justify resting the player. Conversely, when matches are less  
810 valued ( $R = 2$ ), only long injuries justify resting. The decision to rest or play a player critically  
811 depends on the day-specific valuation of matches versus practices.

## 812 **Illustrative Example: Availability Gain Metric for Varying Injury** 813 **Durations**

814 To illustrate the application of our proposed availability gain metric, consider a scenario with 15  
815 observed injuries in the test set, where matches are valued 10 times more than practices ( $R = 10$ ).

### 816 **Scenario 1:**

817 The model correctly predicts one injury, preventing two missed practice days. However, the model  
818 also incorrectly predicts injuries on 50 other occasions, leading to 51 lost practice days (including  
819 the day of true prediction). The availability gain is calculated as:

$$\text{Availability Gain} = \frac{2 - 51}{10} = -4.9 \quad (14)$$

820 This results in a net loss equivalent to 4.9 match days, indicating that the cost of false positives  
821 outweighs the benefit of the correct prediction in this scenario.

### 822 **Scenario 2:**

823 In contrast, if the correctly predicted injury prevents a 30-day absence, including 3 matches and 27  
824 practices, then accounting for the false positives, the net benefit becomes:

$$\text{Availability Gain} = \frac{27 - 51}{10} + 3 = 0.6 \quad (15)$$

825 This represents a net gain in player availability, expressed in match-weighted valuation.

826 These examples highlight how injury duration significantly influences the player availability and how  
827 clubs can adjust  $R$  to reflect specific priorities, such as differentiating between practices, friendly  
828 matches and critical knockout games.

## 829 **Estimating the Survival Curve with the XGBSE Model**

830 The Extreme Gradient Boosting Survival Embedding (XGBSE) model<sup>16</sup> is a three-step procedure  
831 for estimating survival curves from censored survival data. The method combines tree-based model-  
832 ing with logistic regression to provide a robust framework for analyzing survival distributions. The  
833 three steps include: 1. Fitting a censored Accelerated Failure Time model on daily injury observa-  
834 tions. 2. Transforming the feature space based on the fitted Accelerated Failure Time model. 3.  
835 Using logistic regression to estimate binned survival time probabilities in the transformed feature  
836 space to estimate the survival curve.

837 **Step 1: Fitting the Accelerated Failure Time Model**

838 The first step involves fitting the following Accelerated Failure Time model to the data:

$$\ln M_{i,t} = \mathcal{T}(X_{i,t}) + \varepsilon_{i,t},$$

839 where  $M_{i,t}$  represents the number of minutes player  $i$  plays on day  $t$  until injury,  $X_{i,t}$  is the cor-  
 840 responding feature vector,  $\mathcal{T}(X)$  is a tree-ensemble model, and  $\varepsilon_{i,t} \sim \mathcal{N}(0, \sigma^2)$  is Gaussian noise,  
 841 conditionally independent and identically distributed. Injuries are rare events, and for days with-  
 842 out injuries, the actual survival time is unknown. To handle this, a censored likelihood approach is  
 843 employed.

844 The density of an observation, conditioned on the injury indicator  $Y_{i,t}$ , is given by:

$$p(\ln M_{i,t} = m | X_{i,t}, Y_{i,t}; \mathcal{T}) = \phi\left(\frac{m - \mathcal{T}(X_{i,t})}{\sigma}\right)^{Y_{i,t}} \cdot \left[1 - \Phi\left(\frac{m - \mathcal{T}(X_{i,t})}{\sigma}\right)\right]^{1 - Y_{i,t}},$$

845 where  $\phi$  and  $\Phi$  denote the standard normal probability density and cumulative distribution func-  
 846 tions, respectively. The tree ensemble  $\mathcal{T}$  is estimated by minimizing the negative log-likelihood:

$$-\ln \mathcal{L}(\mathcal{T}) = -\sum_{i=1}^N \sum_{t \in \mathcal{I}_i} \ln p(\ln M_{i,t} | X_{i,t}, Y_{i,t}; \mathcal{T}),$$

847 where  $\mathcal{I}_i$  is the set of all observations for individual  $i$ . This approach allows the model to handle  
 848 right-censored data effectively, where exact injury times are unavailable.

849 **Step 2: Feature Space Transformation**

850 After fitting the Accelerated Failure Time model, the tree ensemble  $\mathcal{T}$  is used to transform the  
 851 feature space. Each of the  $K$  trees in the ensemble has  $n_k$  leaf nodes, and each observation is  
 852 assigned to exactly one node per tree. This creates a new high-dimensional, sparse feature space  
 853 represented by indicator variables  $W \in \{0, 1\}^{N \times \sum_{k=1}^K n_k}$ . Each row of  $W$  contains  $K$  indicator  
 854 variables set to one, corresponding to the leaf nodes an observation falls into, while the remaining  
 855 entries are set to zero.

856 **Step 3: Logistic Regression on Binned Survival Times**

857 The transformed feature space  $W$  is then used to fit logistic regressions on binned survival time  
 858 indicators. Observations are grouped into five disjoint time intervals, each spanning 30 minutes,  
 859 denoted as  $[30 \cdot (r - 1), 30 \cdot r)$ . For each interval  $r$ , the logistic regression estimates the conditional  
 860 probability  $q_{i,t}^{(r)}$  of injury occurring within that interval, given that the individual survived until its  
 861 start.

862 The survival function is computed by combining these probabilities:

$$S(\min_{i,t}) = \prod_{r=1}^5 \left(1 - q_{i,t}^{(r)}\right)^{\mathbb{1}_{\{m_{i,t} \in [30 \cdot (r-1), 30 \cdot r)\}}}$$

863 This formulation ensures that the model accurately accounts for survival probabilities over time,  
 864 incorporating both censored and uncensored observations.

## 865 **Augmenting Activity Load Records by Imputing National Team** 866 **Records**

867 To address missing training load data, unreported days were classified as either rest days or missing  
868 values based on predefined rules (see Supplementary Table S2). For players on national duty,  
869 training load was imputed using a three-step process. First, the distance covered in national matches  
870 was predicted using a mixed-effects model that incorporated match minutes and player position.  
871 Second, training loads for national training days were estimated using a gradient boosting model  
872 with predictors such as time since the last match, time until the next match, and player position.  
873 Finally, additional training exposure variables were derived using another gradient boosting model,  
874 which used the imputed distances as inputs and accounted for player-specific variability.

## 875 Supplementary Figures

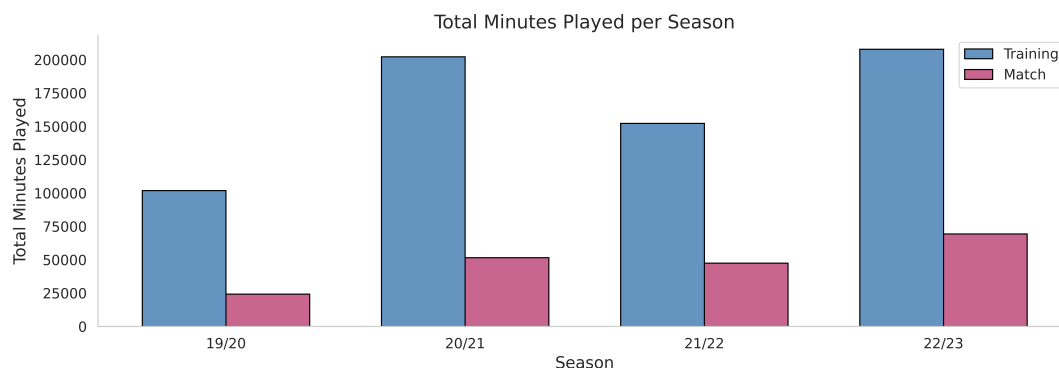


Figure S1: **Played Match and Practice Minutes.** The total time played in minutes per season.

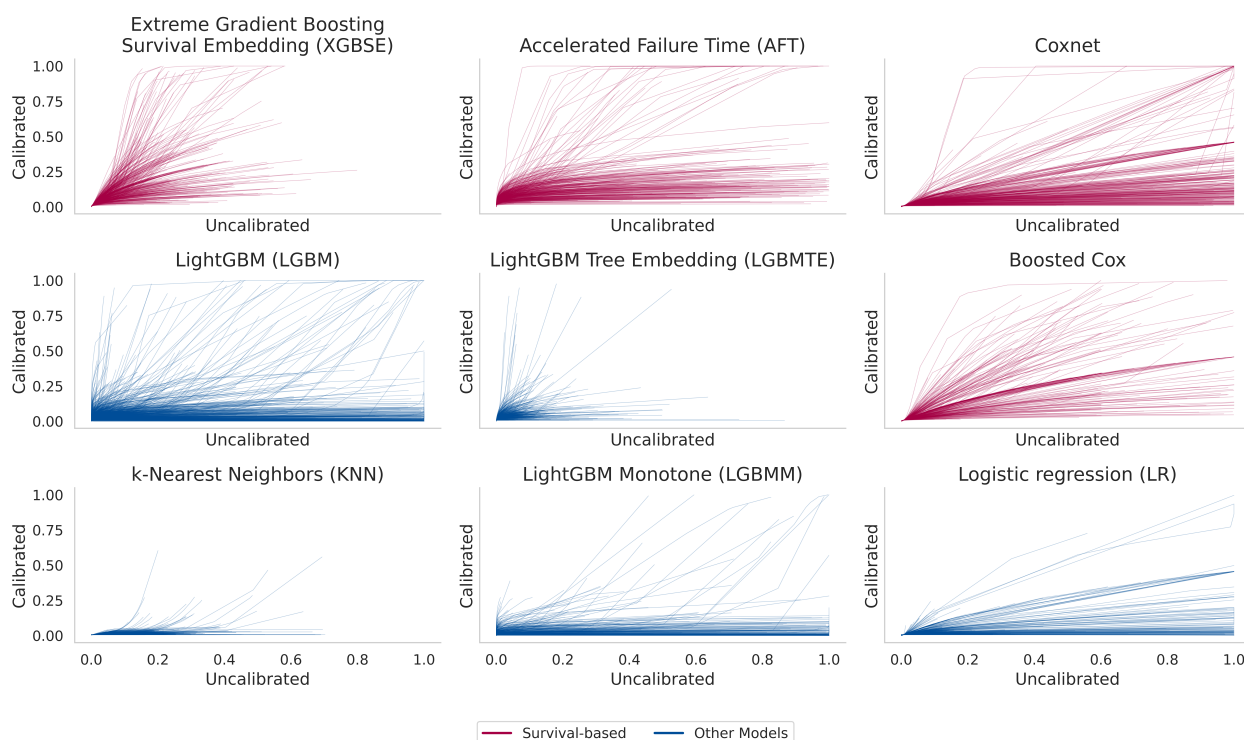


Figure S2: **Beta Calibration Transformations Over Sample Splits.** Each curve represents one fitted beta calibration function on the validation set. One curve is calculated for each of the 400 block splits. Survival-based models are given in red and standard classification models in blue.

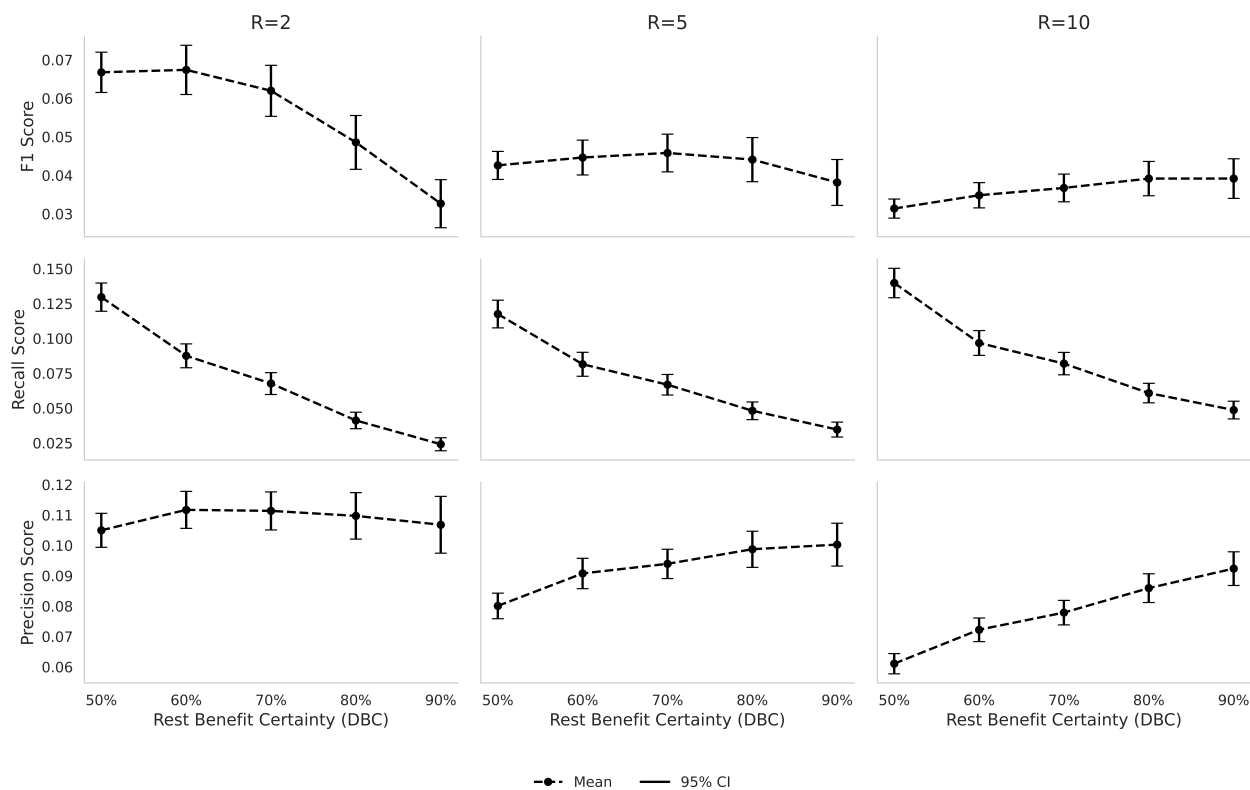


Figure S3: **Mean F1-Score, Precision, and Recall Over Sample Splits.** Each point represents the mean of the respective metric evaluated at the respective Rest Benefit Certainty (RBC) over the 400 sample splits. Confidence intervals represent the 95% confidence intervals around the respective means.

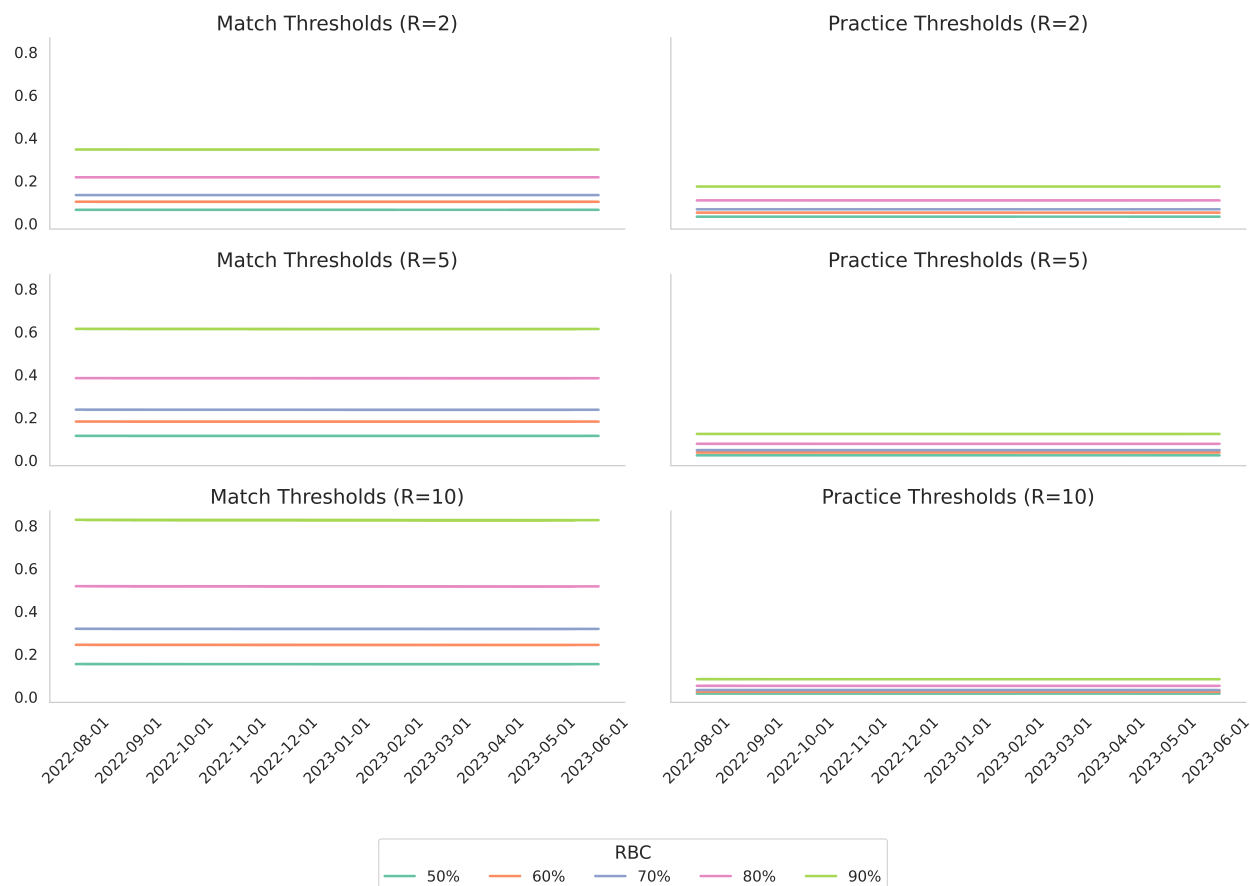


Figure S4: **Decision Thresholds According to Match Valuation and Type of Activity for the 2022/2023 Season.** Each line represents a match or practice threshold, depending on the Rest Benefit Certainty (RBC) and match valuation ( $R$ ) for the 2022/2023 season. By construction of the RBC, the order of the lines is the same across plots, only the magnitudes change.

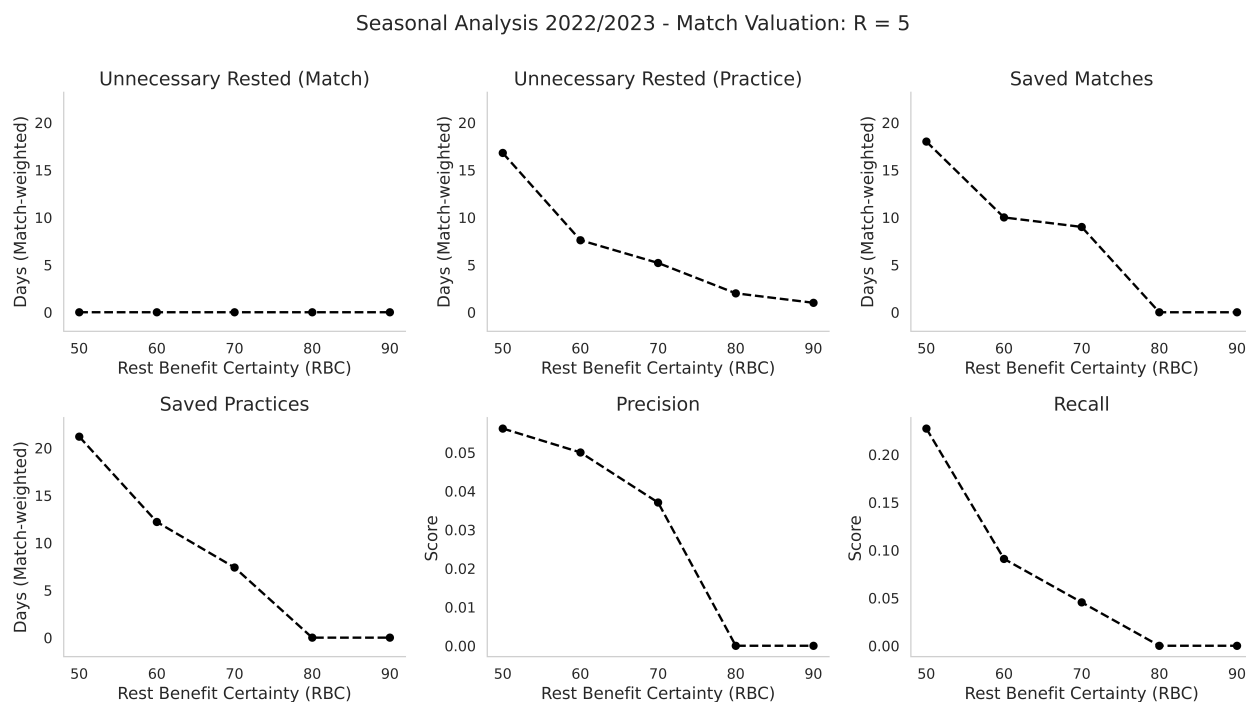


Figure S5: **Additional Evaluation Metrics for Predicting the 2022/2023 Season with R = 5.** Each point represents the respective metric evaluated at the respective Rest Benefit Certainty (RBC) for the 2022/2023 season.

Seasonal Analysis 2022/2023 - Match Valuation:  $R = 2$

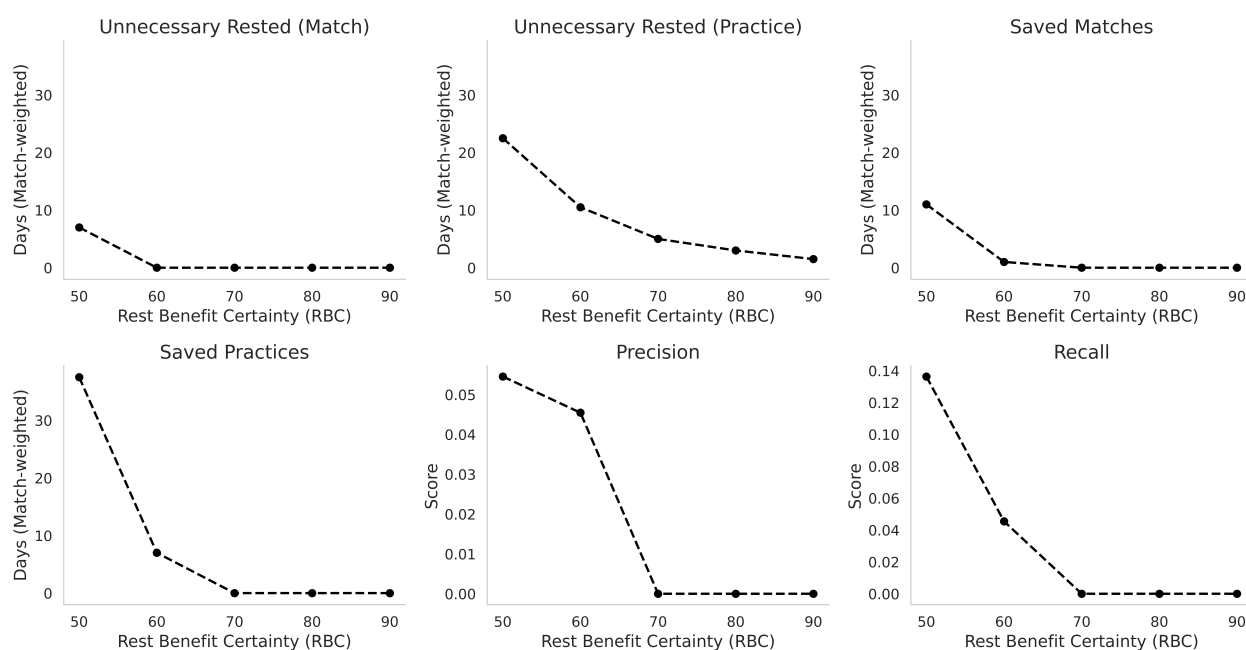


Figure S6: **Additional Evaluation Metrics for Predicting the 2022/2023 Season with  $R = 2$ .** Each point represents the respective metric evaluated at the respective Rest Benefit Certainty (RBC) for the 2022/2023 season.

Seasonal Analysis 2022/2023 - Match Valuation:  $R = 10$

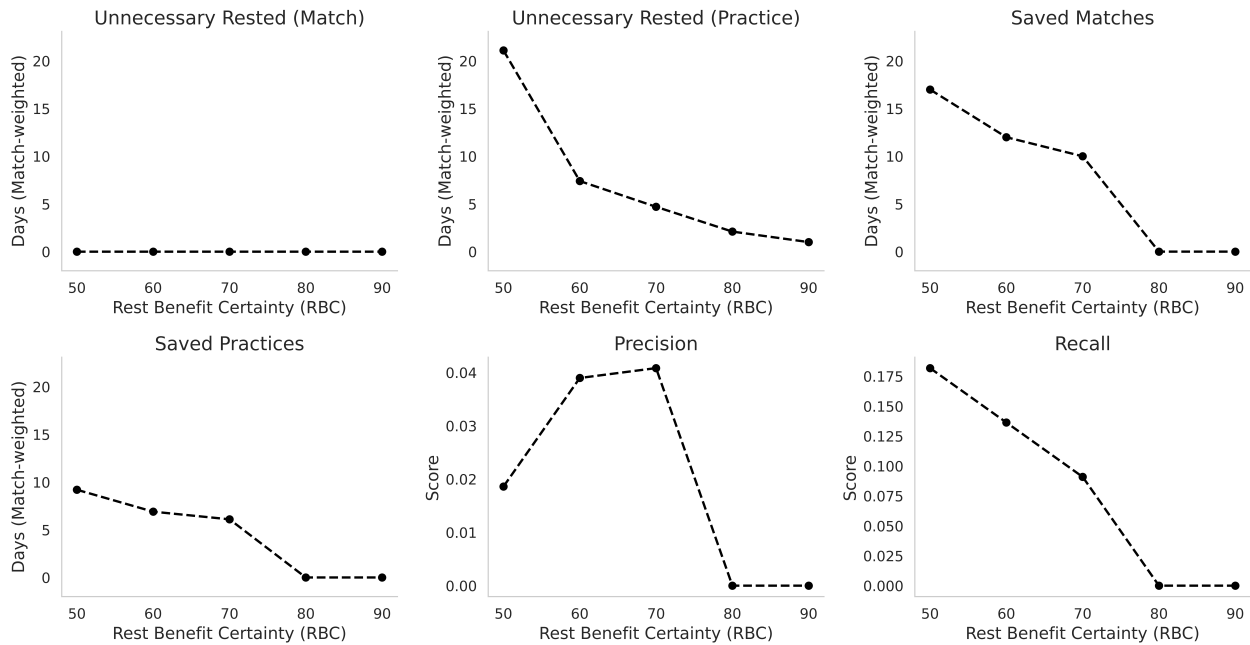


Figure S7: **Additional Evaluation Metrics for Predicting the 2022/2023 Season with  $R = 10$ .** Each point represents the respective metric evaluated at the respective Rest Benefit Certainty (RBC) for the 2022/2023 season.

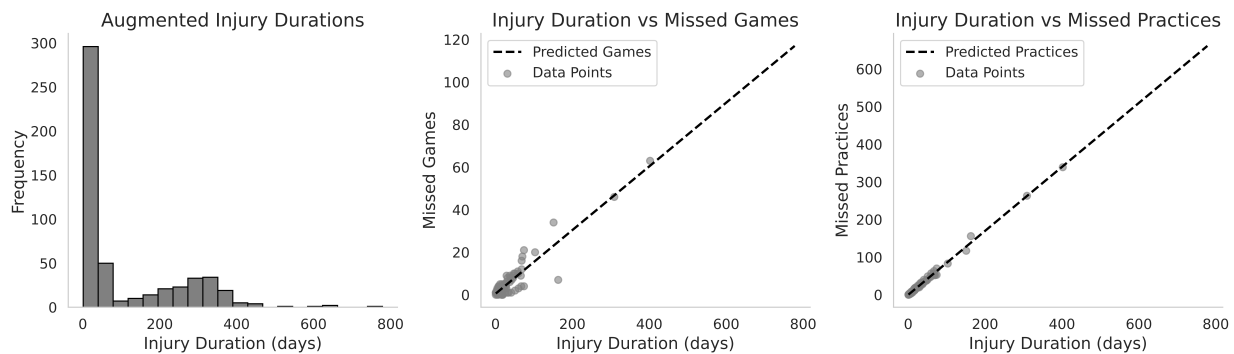


Figure S8: **Augmented Injury Duration Data.** The figure shows the Distribution of the augmented injury data as well as an exemplarily fit of the number of missed games as well as missed practices dependent on the injury duration.

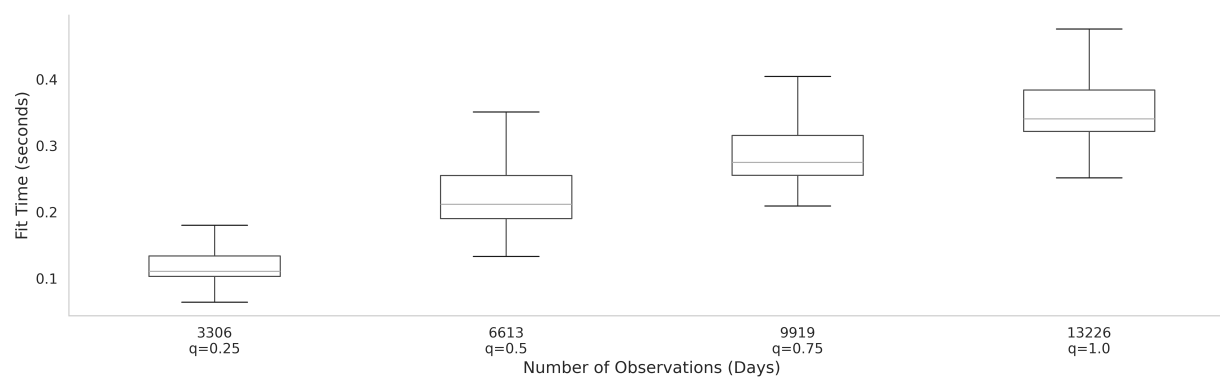


Figure S9: **Fitting Times in Seconds of the XGBSE Model.** The plot shows the fitting times of the XGBSE model on 25%, 50%, 75%, and 100% of the data. Fitting was repeated 400 times. If  $q < 1$ , the data was sampled randomly with stratification for each of the 400 iterations.

## 876 Supplementary Tables

Table S1: **Scraped Non-Contact Injury Types.** The table represents all injuries we used from Soccerdonna<sup>39</sup> in order to augment the injury duration data set as outlined within the Method section.

<b>Term at Soccerdonna</b>	<b>English Translation</b>	<b>Reason for Selection as Non-Contact Injury</b>
Achillessehnenprobleme	Achilles tendon problems	Typically occur without external force.
Achillessehnenreizung	Achilles tendon irritation	Overuse or strain without direct impact.
Achillesfersenprobleme	Achilles heel problems	Result of repetitive stress or strain.
Adduktorenbeschwerden	Adductor issues	Caused by overuse or stretching.
Adduktorenverletzung	Adductor injury	Commonly associated with overexertion.
Außenbandriss Knie	Lateral ligament tear (knee)	Often results from pivoting or twisting.
Bänderriss	Ligament tear	Non-contact due to sudden motion.
Bänderverletzung	Ligament injury	Often caused by excessive strain.
Ermüdungsbruch	Stress fracture	Overuse injury, no external trauma.
Fersneprobleme	Heel problems	Typically from overuse or repetitive stress.
Hüftprobleme	Hip problems	Overload-related without external contact.
Innenbanddehnung Knie	Medial ligament strain (knee)	Twisting or turning movements.
Innenbandriss Knie	Medial ligament tear (knee)	Often from sudden direction change.
Innenbandverletzung	Medial ligament injury	Related to overextension or twisting.

<b>Term at Soccerdonna</b>	<b>English Translation</b>	<b>Reason for Selection as Non-Contact Injury</b>
Kreuzbandriss	ACL tear	Sudden pivot or jump landing.
Leistenprobleme	Groin problems	Typically from overuse or stretching.
Leistenverletzung	Groin injury	Result of excessive strain or motion.
Meniskusschaden	Meniscus damage	Twisting or overloading of the joint.
Muskelfaserriss	Muscle fiber tear	Sudden forceful contraction or stretch.
Muskelfaserriss im Adduktorenbereich	Muscle fiber tear in adductors	Overuse or abrupt movement.
Muskelverletzung	Muscle injury	Common from overuse or strain.
Muskelverhärtung	Muscle hardening	Typically due to overexertion.
Muskelzerrung	Muscle strain	Overstretching of muscle fibers.
Oberschenkelmuskelriss	Thigh muscle tear	Result of high-force contraction.
Oberschenkelprobleme	Thigh problems	Usually overuse or overload.
Oberschenkelzerrung	Thigh strain	Overstretching or repetitive use.
Patellasehnenriss	Patellar tendon tear	Often caused by overloading the knee.
Patellasehnenreizung	Patellar tendonitis	Overuse without direct impact.
Probleme mit Hüftbeuger	Hip flexor problems	Result of overuse or strain.
Schambeinentzündung	Pubic bone inflammation	Overload from repetitive motion.
Seitenband-Verletzung	Lateral ligament injury	Caused by twisting or pivoting.
sehnenentzündung	Tendon inflammation	Overuse without direct trauma.

<b>Term at Soccerdonna</b>	<b>English Translation</b>	<b>Reason for Selection as Non-Contact Injury</b>
Teilschädigung Kreuzband	am Partial ACL damage	Sudden pivot or stop.
Wadenprobleme	Calf problems	Overuse or sudden strain.
Wadenzerrung	Calf strain	Overstretching or fatigue.
Zerrung	Strain	Overextension of tissues.
Zerrung im Oberschenkel- und Gesäßmuskulatur	Thigh and gluteal muscle strain	Overstretching or overload.

Table S2: **Rules for data pre-processing.** The table summarizes how missing training information was handled during feature engineering, as described in the Method section.

	<b>Handling</b>
<b>Physical Activity</b>	
A day at which the total running distance covered is greater than zero at a specific day	Reported value
<b>Rest day</b>	
A day at which the total running distance is not reported, with a match on the previous day	Zero training load
A day at which the total running distance is not reported, flanked by non-zero values on both the previous and next day	Zero training load
2 (3) consecutive days for which the total running distance is not reported and is preceded and succeeded by a reported day	Zero training load
The sequence of days for which the total running distance is not reported including the 24th of December	Zero training load
The sequence of days for which the total running distance is not reported including the 1st of January	Zero training load
The sequence of days for which the total running distance is not reported including the 1st of July 2020 and 2022	Zero training load
<b>National team</b>	

The sequence of days for which the total running distance is not reported but one of the days is associated with a national match	Imputed but only used as information for previous exposure
<b>Injury recovery</b>	
The sequence of days from the day after an injury to 3 weeks after the first day of recorded physical activity after the injury	Dropped
<b>Unknown activity</b>	
Days for which none of the other rules apply to. Includes COVID19 pandemic period.	Dropped

Table S3: **Engineered Features Description.** The table describes the engineered features that were used for model training, as described in the Method section.

Feature Name	Description	Calculation Method
<b>Categorized Training Load</b>		
Distance Covered per Minute	Distance per Minute covered in meters during training or match play.	Directly measured using player tracking data.
High Metabolic Load Distance (HMLD) per Minute	Distance per Minute covered when metabolic power exceeds 25.5 W/kg, typically occurring at speeds above 5.5 m/s.	Directly measured using player tracking data.
Number of Accelerations per Minute	Number of acceleration events per Minute.	Directly measured using player tracking data.
Number of Decelerations per Minute	Number of deceleration events per Minute.	Directly measured using player tracking data.
Distance Covered during Acceleration per Minute	Distance covered per Minute during acceleration events.	Directly measured using player tracking data.
Distance Covered during Deceleration per Minute	Distance covered per Minute during deceleration events.	Directly measured using player tracking data.
High-Speed Running Distance per Minute (Relative HSR)	Distance covered per Minute while running at speeds above 18 km/h.	Directly measured using player tracking data.
<b>Accumulated Training Load</b>		

Continued on next page

**Table S3 – continued from previous page**

Feature Name	Description	Calculation Method
EWMA	Exponentially Weighted Moving Average (EWMA) of past training loads. Provides a smoothed estimate with greater weight on recent loads.	$\alpha e_{i,t} + (1 - \alpha)EWMA_{i,t-1}^j$ with $\alpha = 2/11$ as in <sup>6</sup>
ACWR	Acute Chronic Workload Ratio (ACWR) of the respective training loads. Assesses the balance between acute and chronic load, indicating potential injury risk.	$\frac{\text{Acute Load (last 6 days)}}{\text{Chronic Load (last 28 days)}}$
MWR	Monotonic Workload Ratio (MWR) of the respective training loads. Evaluates training load stability by considering the mean and standard deviation over the last 7 days.	$\frac{\text{Mean of Last 7 Days}}{\text{Standard Deviation of Last 7 Days}}$

Table S4: **Hyperparameter Search Spaces.** This table summarizes the hyperparameter search spaces for each classification model that was used as comparison against the survival-based models. They were applied with Bayesian hyperparameter optimization with 60 iterations for each model. Integer[a,b] indicates all integers values from a to b. Real[a,b] indicates all real values from a to b. Categorical[a,b,c] indicates the values a,b, and c. The option log-uniform indicates that the respective values are drawn uniformly from the logarithm of the bounds.

Model Name	Hyperparameter	Range/Type
<b>k-Nearest Neighbors (KNN)</b>		
Number of Neighbors ( $n\_neighbors$ )	Number of nearest neighbors to consider for classification or regression.	Integer [3, 20]
Weights	Weight function used in prediction.	Categorical ['uniform', 'distance']
Distance Metric ( $p$ )	Distance metric for the tree (Manhattan or Euclidean).	Categorical [1, 2]
<b>LightGBM</b>		
Number of Estimators ( $n\_estimators$ )	Number of boosting iterations.	Integer [20, 300]

Continued on next page

**Table S4 – continued from previous page**

<b>Model Name</b>	<b>Hyperparameter</b>	<b>Range/Type</b>
Maximum Depth ( <i>max_depth</i> )	Maximum depth of the tree.	Integer [3, 8]
Learning Rate ( <i>learning_rate</i> )	Shrinkage parameter for boosting.	Real [0.1, 0.5]
Subsample ( <i>subsample</i> )	Fraction of samples used for fitting.	Real [0.5, 1.0]
Column Subsampling ( <i>colsample_bytree</i> )	Fraction of columns sampled for each tree.	Real [0.5, 1.0]
<b>LightGBM with Monotone Activity Duration</b>		
Number of Estimators ( <i>n_estimators</i> )	Number of boosting iterations.	Integer [20, 200]
Maximum Depth ( <i>max_depth</i> )	Maximum depth of the tree.	Integer [3, 8]
Learning Rate ( <i>learning_rate</i> )	Shrinkage parameter for boosting.	Real [0.1, 0.5]
Subsample ( <i>subsample</i> )	Fraction of samples used for fitting.	Real [0.5, 1.0]
Column Subsampling ( <i>colsample_bytree</i> )	Fraction of columns sampled for each tree.	Real [0.5, 1.0]
<b>Logistic Regression (LR)</b>		
Regularization Parameter ( <i>C</i> )	Inverse of regularization strength.	Real [1e-4, 1e2, log-uniform]
<b>LightGBM with Tree Embedding</b>		
Number of Estimators ( <i>n_estimators</i> )	Number of boosting iterations (LightGBM).	Integer [20, 300]
Maximum Depth ( <i>max_depth</i> )	Maximum depth of the tree (LightGBM).	Integer [3, 8]
Learning Rate ( <i>learning_rate</i> )	Shrinkage parameter for boosting (LightGBM).	Real [0.1, 0.5]
Subsample ( <i>subsample</i> )	Fraction of samples used for fitting (LightGBM).	Real [0.5, 1.0]
Column Subsampling ( <i>colsample_bytree</i> )	Fraction of columns sampled for each tree (LightGBM).	Real [0.5, 1.0]

Continued on next page

**Table S4 – continued from previous page**

<b>Model Name</b>	<b>Hyperparameter</b>	<b>Range/Type</b>
Regularization Parameter ( $C$ )	Inverse of regularization strength (Logistic Regression).	Real [1e-4, 1e2, log-uniform]