

Accepted author manuscript version reprinted, by permission, from International Journal of Sports Physiology and Performance (IJSPP), 2024,
<https://doi.org/10.1123/ijsp.2024-0247>. © Human Kinetics, Inc.

1 **An educational review on machine learning: a SWOT analysis for implementing**
2 **machine learning techniques in football.**

3

4 Marco Beato^{1*}, Mohamed Hisham Jaward², George P. Nassis^{3,4}, Pedro Figueiredo^{3,5}, Filipe
5 Manuel Clemente^{6,7,8}, Peter Krstrup^{4,9}

6

7

8 **Affiliations**

9 1. School of Allied Health Sciences, University of Suffolk, Ipswich, United Kingdom.

10 2. School of School of Technology, Business and Arts, University of Suffolk, Ipswich,
11 United Kingdom.

12 3. Physical Education Department, United Arab Emirates University, Al Ain, United
13 Arab Emirates.

14 4. Department of Sports Science and Clinical Biomechanics, Sport and Health
15 Sciences Cluster (SHSC), University of Southern Denmark, Denmark

16 5. Research Center in Sports Sciences, Health, Sciences and Human Development,
17 CIDESD, Vila Real, Portugal.

18 6. Escola Superior Desporto e Lazer, Instituto Politécnico de Viana do Castelo, Rua
19 Escola Industrial e Comercial de Nun'Álvares, 4900-347 Viana do Castelo,
20 Portugal.

21 7. Gdansk University of Physical Education and Sport, 80-336 Gdańsk, Poland.

22 8. Sport physical activity and health research innovation and technology center
23 (SPRINT), 4900-347 Viana do Castelo, Portugal.

24 9. Danish Institute for Advanced Study (DIAS), University of Southern Denmark,
25 Odense, Denmark.

26

27

28 *Corresponding author

29 Marco Beato, School of Allied Health Sciences, University of Suffolk, Ipswich, United
30 Kingdom, email: m.beato@uos.ac.uk

31

32 **Short title/running head:** A SWOT analysis on machine learning in football.

33

34

35

36 **Abstract**

37 **Purpose:** The abundance of data in football presents both opportunities and challenges for
38 decision-making. Consequently, this review has two primary objectives: first, to provide
39 practitioners with a concise overview of the characteristics of machine learning (ML) analysis;
40 and second, to conduct a Strengths, Weaknesses, Opportunities, and Threats (SWOT) analysis
41 regarding the implementation of ML techniques in professional football clubs. This review
42 explains the difference between artificial intelligence and ML, and the difference between ML
43 and statistical analysis. Moreover, we summarize and explain the characteristics of ML
44 learning approaches such as supervised learning, unsupervised learning and reinforcement
45 learning. Finally, we present an example of SWOT analysis, which suggests some actions to
46 be considered in applying ML techniques by the medical and sport science staff working in
47 football. Specifically, four dimensions were presented namely the use of strengths to create
48 opportunities and make the most of them, the use of strengths to avoid threats, work on
49 weaknesses to take advantage of opportunities, and upgrade weaknesses to avoid threats.

50 **Conclusion:** ML analysis can be an invaluable ally for football clubs, sport science and medical
51 departments due to its ability to analyze vast amounts of data and extract meaningful insights.
52 Moreover, ML can enhance performance **by assessing the risk of injury occurrence,**
53 physiological parameters, physical fitness, and optimizing training, recommending strategies
54 based on opponent analysis, and identifying talent and assessing player suitability.

55

56 **Key Points:** Strengths, Weaknesses, Opportunities, Threats, decision-making, performance
57 prediction, injury risk assessment, Soccer

58

59

60

61

62

63

64

65

66

67 **INTRODUCTION**

68 The decision-making process plays a critical role for practitioners working in football.
69 Practitioners aim to optimize the training process, testing protocols, physiological parameters,

70 physical readiness, and match strategies to increase the probability of success.¹⁻³ In the last
71 decades, technology has allowed sports scientists and performance analysts to collect larger
72 volumes of data compared to the past,⁴⁻⁶ and use them in conjunction with their experience
73 and the most relevant scientific evidence to make informed decisions. These data have been
74 typically analyzed using visualizations and statistical methods. Nevertheless, challenges arise
75 when determining how to effectively select variables and handle larger datasets derived from
76 multiple sources and instruments. In recent years, artificial intelligence (AI) and machine
77 learning (ML) have become more pervasive in football⁶⁻⁸. Although the use of AI and ML are
78 common in our contemporary society, some confusion exists between the two terms. AI can be
79 briefly defined as “the theory and development of computer systems able to perform tasks
80 normally requiring human intelligence”,⁹ while ML refers to “the technologies and algorithms
81 that enable systems to identify patterns, make decisions, and improve themselves through
82 experience (training)”¹⁰ and it is a subset of AI. ML can find several applications in football,
83 for example, to facilitate decision-making, performance prediction, technical and tactical
84 pattern recognition, game activity/analytics, talent identification, and injury risk
85 assessment.^{7,11,12}

86
87 Data mining is the process of sorting through large data sets to identify patterns and
88 relationships.¹³ Through data mining and ML (which focuses on creating algorithms that can
89 learn and predict from given data)⁷ football practitioners (e.g., sports scientists and coaches)
90 can make informed decisions to enhance physiological parameters, physical development,
91 reduce fatigue, increase readiness, and match performance. A recent review reported that ML
92 can be used to determine the parameters that affect (i.e., explainability, which means that a
93 model and its output can be explained and make sense to a human being) wellness and fitness,
94 which can be later on manipulated by football practitioners.⁷ ML regression can determine the
95 contribution of players’ anthropometric characteristics to physical performance, such as
96 sprinting and aerobic fitness.¹⁴ Furthermore, ML can be used to assess the relationship between
97 well-being parameters and training load and match performance. However, it showed a limited
98 predictive capacity of such parameters to determine internal and external load.¹⁵ ML analysis
99 can be used for determining technical and tactical outcomes, for instance, to analyze the team
100 pattern or the effectiveness of passing strategies.⁷ ML was used to estimate players’ passing
101 skills to make predictions for the following season,¹⁶ which coaches and performance analysts
102 could use for scouting objectives. Moreover, multiple ML algorithms were used by Jamil et
103 al.,¹⁷ to classify elite and sub-elite goalkeepers (GK) in professional men's football, suggesting

104 that a GK's ability with their feet and not necessarily their hands are what distinguishes the elite
105 GK's from the sub-elite. Another area in which ML can be used is talent identification, which
106 is one of the more critical challenges for football clubs. In this specific context, technical and
107 tactical variables, together with psychological and physical variables can be assessed to
108 determine the talent predictors that coaches need to monitor and develop.^{18,19} Such information
109 may impact the productivity (in terms of talent) of football academies and related clubs.
110 Certainly, ML holds the promise to overcome the constraints of conventional reductionism
111 approaches, enabling the concurrent integration of diverse data sources. It may play a pivotal
112 role in gathering a comprehensive understanding of the game by bridging gaps across physical,
113 physiological, technical, and tactical dimensions, while simultaneously contextualizing the
114 information and actively pursuing integrative models. This advanced approach may accelerate
115 analyses but also potentially heightens accuracy, thereby strengthening decision-making
116 processes in coaching, player development, and overall team performance.

117

118 Research in the field of ML for identifying injury risks and associated factors has been steadily
119 growing over the years, as evidenced by a recent systematic review¹¹. For instance, in a study
120 by Oliver et al.,²⁰ involving 355 elite youth football players, decision tree algorithms displayed
121 an overall accuracy that was not significantly superior to statistical logistic regression in
122 detecting injuries. However, ML (e.g., decision tree) demonstrated increased sensitivity in this
123 context. In contrast, a study by Rommers et al.,²¹ which employed extreme gradient boosting
124 algorithms on a larger sample of 734 youth players, revealed promising results. The ML
125 algorithm successfully identified injured players in the hold-out test sample with 85%
126 precision, 85% recall (sensitivity), and 85% accuracy.²¹ Additionally, the same study²¹
127 achieved reasonably high accuracy in distinguishing between overuse and acute injuries based
128 on pre-season measures. Hence, beyond predicting potential injuries, ML has the potential to
129 categorize them effectively. This capability provides additional insights for rapidly
130 constructing models in subsequent stages of interpretation. Furthermore, it facilitates
131 interaction with potential injury mechanisms and factors that may influence the overall risk.²²

132

133 To successfully implement ML in football, practitioners need to address the integration into
134 medical, sport science, and coaching departments. A strategic management plan, anchored by
135 a Strengths, Weaknesses, Opportunities, and Threats (SWOT) analysis, is vital to evaluate the
136 team's internal capabilities and external possibilities concerning ML. This analysis will inform
137 strategic decisions, leveraging strengths to harness opportunities or neutralize threats, and

138 improving weaknesses to support ML adoption. A team's preparedness to adopt ML is crucial,
139 as it can significantly refine their strategic approach to ML utilization, ensuring a more
140 effective and efficient integration. This condensed strategy enables teams to navigate the
141 complexities of ML implementation in the competitive sports environment.

142

143 In the dynamic field of football, the profusion of data creates a spectrum of possibilities and
144 hurdles in the decision-making process. Addressing this, the review unfolds in two distinct
145 parts: the first segment offers practitioners a streamlined synopsis of ML analysis features; the
146 second segment presents a comprehensive SWOT analysis, assessing the practicality and
147 impact of integrating ML methodologies within the ecosystem of professional football clubs.

148

149 **MACHINE LEARNING**

150 **Difference between AI and ML**

151 Before delving into ML in football, it is important to appreciate the evolution of ML to the
152 modern form used to solve many real-world problems. As mentioned in the introduction, ML
153 constitutes a subset within the broader field of AI. Modern AI gained prominence in the early
154 1940s and the seminal work of McCulloch and Pitts is considered as the first work on the
155 artificial neuron (they defined a mathematical computation model similar to neural networks).²³
156 Various AI initiatives aim to emulate human intelligence through computational models based
157 on artificial neurons. Consequently, AI encompasses a wide spectrum of tasks and issues, in
158 contrast to ML, where the primary objective is the development of algorithms tailored for
159 specific tasks. Frequently, everyday tasks can be formulated as either regression or
160 classification problems, and ML endeavors to address these challenges systematically.

161

162 **Difference between ML and statistical analysis**

163 In numerous data science scenarios, the principal goals are inference and prediction. Inference
164 involves creating a mathematical model of the data-generation process to formalize
165 understanding or test hypotheses regarding system behavior. As an example in football, Zeki
166 et al.²⁴ infer the neuromuscular fatigue imposed on players after a football match based on
167 measurements such as the players' heart rate, accelerations, and distance traveled. Prediction
168 aims at forecasting unobserved outcomes or future behavior, such as whether a football player
169 will likely develop an injury in a future game. In a typical research project or applied setting,
170 both inference and prediction can be of value - we want to know how the system works and
171 what will happen next.²⁵

172

173 Many methods from statistics and ML may, in principle, be used for both prediction and
174 inference. However, statistical methods have a long-standing focus on inference, which is
175 achieved through the creation and fitting of a probabilistic model onto the data.²⁶ The model
176 allows us to compute a quantitative measure of confidence that a discovered relationship
177 describes a 'true' effect, and so unlikely to result from noise or disturbances. In contrast, ML
178 emphasizes prediction, employing learning algorithms to identify patterns in complex and big
179 datasets.²⁶ ML techniques prove particularly advantageous when dealing with situations where
180 the number of input variables surpasses the number of samples, as opposed to scenarios with
181 more samples than input variables. ML operates with minimal assumptions about data-
182 generating systems, exhibiting efficacy even in instances where data collection lacks a
183 meticulously controlled experimental design or involves intricate nonlinear interactions. Also,
184 ML allows for the interdependence of data points and facilitates the identification of hidden
185 targets/groups without needing a subjective setting while providing an error estimation.²⁷
186 However, despite achieving compelling predictive outcomes, the limited interpretability of
187 numerous ML solutions poses challenges in directly addressing specific problems and applying
188 them in safety-critical applications. Often, statistical methods, including hypothesis testing, are
189 employed to validate ML outcomes, and the relative performance of ML methods is commonly
190 compared using hypothesis testing approaches.

191

192 Classical statistics and ML diverge in terms of computational tractability as the number of
193 variables per subject increases.²⁶ Classical statistical modeling, originally designed for datasets
194 with a limited number of input variables and sample sizes considered small to moderate by
195 contemporary standards, encounters challenges as the complexity of the relationships among
196 numerous input variables increases. Consequently, statistical inferences become less precise
197 and the boundary between statistical and ML approaches becomes hazier.

198

199 **SUPERVISED AND UNSUPERVISED ML METHODS**

200 **Supervised learning**

201 In supervised learning, a model is derived from a dataset that incorporates features and labels,
202 with both entities employed during the training phase (see Table 1).^{17,28,29} Once the model is
203 trained, it predicts the label corresponding to the input features (values that a supervised model
204 uses) when presented with unseen input (the value we want the model to predict). The
205 supervisor oversees the learner's every move, dictating precise actions for every situation until

206 the learner masters the mapping from situations to actions. While working under such close
207 supervision may seem restrictive, the process is relatively straightforward – quickly
208 recognizing patterns and replicating the supervisor's actions ensures compliance.

209

210 In supervised ML, the supervisory aspect is crucial, as it forces the model to learn parameters
211 of the model such that the output given by the model is close to the desired output indicated by
212 the label. In probabilistic terms, the focus is typically on estimating the conditional probability
213 of a label given specific input features. While supervised learning represents just one paradigm
214 among several, it predominates in the success of ML applications across various domains.^{30,31}
215 This prevalence is attributed, in part, to the fact that many pivotal tasks, such as those listed
216 below, revolve around estimating the probability of an unknown attribute given a specific set
217 of available data:

- 218 • Assess injury risk in elite youth football players using ML.²¹
- 219 • Classifying elite and sub-elite goalkeepers in professional men's football.¹⁷
- 220 • Effective injury forecasting in soccer with GPS training data and ML.³²
- 221 • Predicting the stock price (e.g., of a Club) for the next month based on this month's
222 financial reporting data.²⁹

223 Despite all supervised learning problems being encapsulated by the overarching description
224 of "predicting labels given input features," the methodology assumes diverse forms and
225 necessitates numerous modeling decisions. These decisions hinge on considerations such as
226 the type, size, and quantity of inputs and outputs, leading to the utilization of different models
227 tailored for processing sequences of varying lengths and fixed-length vector representations,
228 among other factors.

229 *****Table 1 here*****

230

231 **Regression and classification**

232 Perhaps the simplest supervised learning task is regression. A typical illustration of a regression
233 problem involves predicting a player transfer market value based on various factors such as
234 age, performance statistics, experience, etc. Goddard (2005) applied regression techniques to
235 forecast goals scored and conceded,³³ leveraging a 25-year dataset on English league football
236 match outcomes. The defining characteristic of a regression problem lies in the form of the
237 target variable. When labels assume arbitrary numerical values, even within a specific interval,

238 the problem is classified as a regression problem. The primary objective is to develop a model
239 that produces predictions closely aligned with the actual numerical label values.

240 In contrast to regression, the output of a classifier takes only a finite number of values. In
241 classification tasks, the model predicts the category (often termed a class) to which a given
242 example belongs from a discrete set of options. For instance, automatic activity classification
243 in sports, like jumping or running. The most basic form of classification is binary classification,
244 where the scenario involves only two classes. While regression employs a regressor to output
245 a numerical value, classification seeks a classifier whose output predicts the assigned class.
246 Despite classification and regression being distinct problems, analogous models are employed
247 to address both sets of challenges. In classification, classes are distinguished using a decision
248 boundary, whereas in regression, efforts are directed towards minimizing the difference
249 between training samples and the values predicted by the boundary.

250

251 **Decision tree**

252 Decision trees stand out as a widely adopted ML technique employed to establish connections
253 between input variables, depicted within the branches and nodes of the tree, and an output value
254 encapsulated in the leaves of the tree. The decision tree is one of the oldest and most popular
255 techniques for supervised learning, which has been developed independently in the statistical³⁴
256 and ML³⁵ communities. These trees find applications in both classification problems, where
257 they produce a category label, and regression problems, where they yield a real number as
258 output. Various algorithms, including the well-established CART (Classification and
259 Regression Tree, which produces only binary Trees) or ID3 (Iterative Dichotomiser 3, which
260 produces decision trees with nodes having more than two children), are employed for fitting
261 decision trees, employing a combination of greedy searching and pruning strategies to ensure
262 the tree effectively fits the training data while also generalizing well to unseen input/output
263 pairs.

264

265 A notable advantage of decision trees lies in their scalability with additional data, resilience to
266 irrelevant features, and interpretability. The choices made at each node facilitate an
267 understanding of the impact of each predictor variable on the ultimate outcome. Random
268 forests operate by constructing a multitude of decision trees during training, utilizing different
269 subsets of the dataset as the training set for each tree.³⁶ In classification scenarios, the final

270 output is determined by the mode of the outputs of each decision tree, while for regression
271 problems, the mean is computed. This approach yields a model with significantly enhanced
272 performance compared to a single decision tree, attributed to reduced overfitting. Nevertheless,
273 the interpretability of the model diminishes, as the decisions at the nodes of the individual trees
274 differ.

275

276 **Support vector machines (SVMs)**

277 Support vector machines (SVMs) are ML models for classification and regression tasks³⁷. In
278 SVM models, the training data is represented as points in space, aiming to delineate distinct
279 categories by a hyperplane (a crucial deciding boundary that partitions the input space into two
280 or more sections) situated as far as possible from the nearest data points. New input instances
281 are subjected to the same mapping as the training data, enabling their categorization based on
282 their position relative to the hyperplane. In instances where the data lacks linear separability,
283 the kernel trick is used. This is a technique employed in SVMs to transform data that is not
284 linearly separable into a higher-dimensional feature space, where it can potentially be separated
285 linearly.³⁸ Extending beyond classification, SVMs can effectively address regression problems
286 by relying on a subset of the training data to formulate regression predictions and is commonly
287 known as support vector regression. Advantages of using SVMs include that they are effective
288 in high dimensional spaces, that they are memory efficient thanks to the use of a subset of
289 training points in the decision function, and finally that they are versatile through the use of
290 different possible kernel functions. On the other hand, using SVMs can have some
291 disadvantages: they do not directly provide probability estimates for classification problems,
292 and correctly optimizing the kernel function and regularization term is essential to avoid
293 overfitting.

294

295 **Neural networks**

296 Neural networks, also known as artificial neural networks, are systems based on a collection
297 of nodes (neurons) designed to algorithmically emulate the interconnections between neurons
298 in the human brain.³⁹ Each neuron can receive signals from other neurons and transmit them to
299 additional neurons, establishing a network of interconnections. The relationship between two
300 neurons is facilitated by an edge or *arrow* (which, represents the weights and biases of linear
301 transformations between the layers), characterized by a weight that signifies the significance

302 of the input from one neuron to the output of the other. Typically, a neural networks comprises
303 an input layer, featuring one neuron per input variable for the model, an output layer with a
304 single neuron providing the classification or regression outcome, and several hidden layers
305 positioned between the input and output layers, each containing a variable number of neurons.
306 An example of the use of neural networks in team sports can be found here, Ruddy et al.,
307 developed predictive modelling of hamstring strain injuries in elite Australian footballers.⁴⁰

308

309 The advantages of using neural networks as classification or regression models are that they
310 usually achieve higher predictive accuracy than other techniques. However, their effectiveness
311 is contingent upon a substantial volume of training data to optimize the model. Furthermore,
312 neural networks lack a guarantee of convergence to a singular solution, rendering them non-
313 deterministic. Importantly, neural networks lack interpretability due to the complexity
314 introduced by numerous layers and neurons, making it challenging to discern the direction and
315 magnitude of the association between each input variable and the output variable through the
316 different weights.

317

318 **Unsupervised learning**

319 The previous sections focused on supervised learning, where a large dataset containing both
320 features and corresponding label values is provided to the model. In this scenario, the
321 supervised learner operates under the guidance of a highly specialized supervisor. In contrast,
322 envisioning the opposite scenario involves working for a supervisor with ambiguous
323 expectations. In this context, the supervisor might furnish a vast dataset and instruct the data
324 scientist to perform some ML algorithms without providing specific guidance. This ambiguity
325 characterizes a class of problems known as unsupervised learning, wherein the range of
326 questions one can pose is limited only by one's creativity. One common question addressed is
327 to find a small number of prototypes that accurately summarize the data (e.g., given a set of
328 players' characteristics, we can group them into categories). This action is typically known as
329 clustering. Another important and exciting recent development in unsupervised learning is the
330 advent of deep generative models. These models aim to estimate the data density, either
331 through explicit or implicit methods.^{41,42}

332

333 **Clustering**

334 Cluster analysis (predictive or descriptive) is an approach that organizes data objects based

335 solely on information inherent in the data describing these objects and their interrelationships.⁴³
336 The primary objective is to assemble objects within a group that exhibit similarity or
337 relatedness while maintaining dissimilarity or unrelatedness to objects in other groups. The
338 efficacy of clustering is contingent upon achieving homogeneity within a group and
339 maximizing dissimilarity between groups, thereby enhancing the distinctiveness of the
340 clustering outcomes. Cluster analysis shares commonalities with other techniques employed
341 for partitioning data objects into groups. It can be perceived as a variant of classification, as it
342 involves labeling objects with class (cluster) labels derived exclusively from the data. In
343 contrast, classification is a supervised process, where new, unlabeled objects receive class
344 labels using a model developed from objects with known class labels. Consequently, cluster
345 analysis is considered a form of unsupervised classification. In ML, the unqualified term
346 "classification" typically refers to the supervised classification discussed in previous sections.

347 There are many types of clustering techniques, but the most common approach is known as K-
348 means. K-means is a prototype-based, partitional clustering technique striving to identify a
349 user-specified number of clusters (K) represented by their centroids. Agglomerative
350 Hierarchical Clustering encompasses a group of closely related techniques that yield a
351 hierarchical clustering. It initiates by treating each point as a singleton cluster and iteratively
352 merges the two closest clusters until a single, overarching cluster remains. Some of these
353 techniques have a natural interpretation in terms of graph-based clustering, while others have
354 an interpretation in terms of a prototype-based approach.

355

356 **Reinforcement Learning**

357 In methodologies of learning discussed in previous sections, predictions were made on models
358 trained with data from a similar distribution, leading to prediction failures when the system
359 underwent significant changes compared to its training state.¹² In such dynamic situations, we
360 could develop an agent that interacts with an environment and takes actions, then our learning
361 paradigm is known as reinforcement learning. This approach finds applications in diverse
362 domains such as evaluating players performance and training,⁴⁴ including robotics, and the
363 development of AI for video. In the recent past, deep reinforcement learning, which applies
364 deep learning to reinforcement learning problems, has surged in popularity. Although not
365 related to football but sport in general, notable works include the groundbreaking deep Q-
366 network, which outperformed humans in Atari games using only visual input,⁴⁵ and the

367 AlphaGo program, which triumphed over the world champion in the board game Go.⁴⁶
368 Reinforcement learning gives a very general statement of a problem in which an agent interacts
369 with an environment over a series of time steps. At each time step, the agent receives some
370 observation from the environment and must choose an action that is subsequently transmitted
371 back to the environment via some mechanism. After each iteration, the agent receives a reward
372 from the environment. The agent then receives a subsequent observation, and chooses a
373 subsequent action, and so on. The behavior of a reinforcement learning agent is governed by a
374 policy. In brief, a policy is just a function that maps from observations of the environment to
375 actions. The goal of reinforcement learning is to produce good policies.

376

377 **Strengths, Weaknesses, Opportunities and Threats (SWOT) analysis**

378 Earlier sections have explored the integration of ML in football, detailing and clarifying the
379 distinct features of ML learning strategies, including supervised, unsupervised, and
380 reinforcement learning. The forthcoming section introduces a SWOT analysis, proposing
381 several considerations for the implementation of ML tactics by football's medical and sports
382 science departments. It specifically outlines four strategic aspects: 1) use strengths to create
383 opportunities and make the most of them, 2) use strengths to avoid threats, 3) work on
384 weaknesses to take advantage of opportunities, 4) upgrade weaknesses to avoid threats. The
385 SWOT analysis process is a valuable tool for organizations and businesses (i.e., clubs) to assess
386 their internal and external environment. Table 2 reports some key needs for conducting a
387 SWOT analysis.

388

389 ***** Table 2 here*****

390

391 ***Practical tips to run a SWOT analysis in football aiming to apply ML***

392 Before applying ML in the team, medical and sport science staff are advised to build a strategic
393 management plan. As part of this plan, they should perform an environmental analysis, which
394 includes scanning the internal and external factors.⁴⁷ The internal factors include analyzing the
395 strengths and weaknesses of their team/organization.⁴⁷ The external factors analysis includes
396 the factors outside the team/organization, the opportunities, and threats of using ML. This is
397 called SWOT analysis and is being used in other domains too.⁴⁷ An example of a SWOT
398 analysis for a top-level football club is presented in Figure 1. We have assumed that the club's
399 top management has adopted ML to improve their senior squad's injury risk assessment

400 strategy. This new approach may bring value provided the team is ready to take advantage of
401 that opportunity (see, Figure 1).

402

403 *****Please, add here Figure 1*****

404

405

406 With regards to the SWOT analysis presented above, we are suggesting some actions to be
407 considered by the medical and sport science staff working in the club. In particular and with
408 regards to the:

409 • Strategic dimension 1. Use strengths to take advantage of opportunities: the supporting
410 team staff can work with top management to convince the coaches of the competitive
411 advantages this new approach may bring to the team. The highly skilled ML staff can
412 work effectively on optimizing systems and building algorithms for injury risk
413 assessment.¹¹

414 • Strategic dimension 2. Use strength to avoid threats: the supports team staff may work
415 on knowledge transfer to the coaches. Simultaneously, the support team staff should
416 receive further education on technical and tactical aspects of football to better
417 understand the game. This will help in accounting for the context when analyzing big
418 data. In turn, this will facilitate the communication of the support team staff with the
419 coaches and optimize knowledge implementation.⁴⁸

420 • Strategic dimension 3. Upgrade weaknesses to take advantage of opportunities:
421 implement a holistic player-centric monitoring system and consider the complexity of
422 injury occurrence.^{8,49} This will help in better interpreting the algorithms.⁵⁰

423 • Strategic dimension 4. Update weaknesses to avoid threats: optimize player's
424 monitoring and integration of ML tools with the existing systems and workflows, while
425 working on knowledge transfer to the coaches.⁵⁰ Build a "bright spot" that will add a
426 competitive advantage to the team.

427

428 *Limitations and future directions*

429 The implementation of ML is not without limitations or barriers. First, ML models require large
430 amounts of high-quality data for training. In football, obtaining comprehensive and accurate
431 data can be challenging due to variations in data collection methods, inconsistencies, and
432 missing information. For instance, limited historical data for specific events (e.g., injuries,

433 specific player movements) can hinder model performance. Second, ML techniques are not
434 guaranteed to provide correct information (e.g., poor model performance, incorrect prediction
435 and therefore, do not always enhance decision-making). Third, many ML algorithms operate
436 as black boxes (if practitioners do not have a specific background in ML), making it difficult
437 to understand how they arrive at specific decisions. In football, coaches and analysts need
438 interpretable models to make informed decisions. Fourth, creating relevant features (input
439 variables) for football-specific tasks can be complex. Deciding which player attributes, team
440 statistics, or match context to include requires domain knowledge. Moreover, football events
441 (e.g., goals, fouls, yellow cards) occur infrequently compared to non-events (e.g., passes, ball
442 possession). This class imbalance affects model training and evaluation. Therefore, techniques
443 like oversampling, undersampling, or using weighted loss functions are necessary to address
444 this issue. Finally, football is highly context-dependent. Player actions depend on the game
445 situation, opponent, field position, and time remaining. ML models must account for these
446 dynamic factors.

447

448 **Practical applications**

449 ML models can analyze player data (such as physical condition, physiological parameters,
450 match performance, and training load) to **assess** the likelihood of injuries. Clubs can use this
451 information to manage player load, optimize recovery, and reduce injury risks. ML algorithms
452 can assess player form by analyzing historical performance data. Clubs can identify players
453 who are in peak form and make informed decisions about team selection. For scouting, ML
454 can analyze player statistics, playing style, and potential fit with the team's tactics. It helps
455 clubs discover talented players and make strategic signings. ML techniques can analyze
456 opponents' playing styles, strengths, and weaknesses. Clubs can use this information to tailor
457 their game plans, identify vulnerabilities, and exploit opponent weaknesses during matches.
458 ML algorithms can evaluate youth players' performance metrics and potential. Clubs can
459 identify promising talents early, nurture their development, and integrate them into the senior
460 team. Finally, clubs that want to build a strategic management plan can use the four dimensions
461 presented in our SWOT analysis such as the use of strengths to create opportunities and make
462 the most of them, the use of strengths to avoid threats, work on weaknesses to take advantage
463 of opportunities, and upgrade weaknesses to avoid threats.

464

465 **Conclusion**

466 This education review provides practitioners with a concise overview of the characteristics of

467 ML analysis and a guide for how to conduct a SWOT analysis regarding the implementation
468 of ML techniques in professional football clubs. This review explains the difference between
469 AI and ML, and the difference between ML and statistical analysis. Furthermore, we explained
470 the characteristics of ML approaches such as supervised learning, unsupervised learning, and
471 reinforcement learning. Finally, we presented an example of a SWOT analysis, which
472 suggested some actions to consider when ML is implemented by medical and sport science
473 staff in football. In conclusion, ML analysis can be an invaluable ally of football clubs and
474 sport science and medical departments due to its ability to analyze vast amounts of data and
475 extract meaningful insights.

476

477 **References**

- 478 1. Jeffries, A. C. *et al.* Development of a revised conceptual framework of physical
479 training for use in research and practice. *Sport. Med.* **52**, 709–724 (2022).
- 480 2. Dello Iacono, A., Beato, M., Unnithan, V. B. & Shushan, T. Programming high-speed
481 and sprint running exposure in football: beliefs and practices of more than 100
482 practitioners worldwide. *Int. J. Sports Physiol. Perform.* 1–16 (2023)
483 doi:10.1123/ijsp.2023-0013.
- 484 3. Beato, M. *et al.* Rationale and practical recommendations for testing protocols in
485 female soccer: a Narrative Review. *J. Strength Cond. Res.* (2023).
- 486 4. Scott, T. U., Scott, T. J. & Kelly, V. G. The validity and reliability of global
487 positioning system in team sport: a brief review. *J. Strength Cond. Res.* **30**, 1470–1490
488 (2016).
- 489 5. Drust, B. & Green, M. Science and football: evaluating the influence of science on
490 performance. *J. Sports Sci.* **31**, 1377–1382 (2013).
- 491 6. Rein, R. & Memmert, D. Big data and tactical analysis in elite soccer: future
492 challenges and opportunities for sports science. *Springerplus* **5**, 1410 (2016).
- 493 7. Rico-González, M., Pino-Ortega, J., Méndez, A., Clemente, F. & Baca, A. Machine
494 learning application in soccer: a systematic review. *Biol. Sport* **40**, 249–263 (2023).
- 495 8. Claudino, J. G. *et al.* Current approaches to the use of artificial intelligence for injury
496 risk assessment and performance prediction in team sports: a systematic review. *Sport.*
497 *Med. - Open* **5**, 28 (2019).
- 498 9. Roth, E. M. & Woods, D. D. Cognitive task analysis: an approach to knowledge
499 acquisition for intelligent system design. in 233–264 (1989). doi:10.1016/B978-0-444-
500 87321-7.50014-4.

- 501 10. Bhatt, G. D. & Zaveri, J. The enabling role of decision support systems in
502 organizational learning. *Decis. Support Syst.* **32**, 297–309 (2002).
- 503 11. Nassis, G., Verhagen, E., Brito, J., Figueiredo, P. & Krstrup, P. A review of machine
504 learning applications in soccer with an emphasis on injury risk. *Biol. Sport* **40**, 233–
505 239 (2023).
- 506 12. Chmait, N. & Westerbeek, H. Artificial Intelligence and Machine Learning in Sport
507 Research: An Introduction for Non-data Scientists. *Front. Sport. Act. Living* **3**, (2021).
- 508 13. Agarwal, S. Data mining: data mining concepts and techniques. in *2013 International*
509 *Conference on Machine Intelligence and Research Advancement* 203–207 (IEEE,
510 2013). doi:10.1109/ICMIRA.2013.45.
- 511 14. Bongiovanni, T. *et al.* Importance of anthropometric features to predict physical
512 performance in elite youth soccer: a machine learning approach. *Res. Sport. Med.* **29**,
513 213–224 (2021).
- 514 15. Campbell, P. G. *et al.* Analysing the predictive capacity and dose-response of wellness
515 in load monitoring. *J. Sports Sci.* **39**, 1339–1347 (2021).
- 516 16. Szczepański, T. & McHale, I. Beyond completion rate: evaluating the passing ability
517 of footballers. *J. R. Stat. Soc.* **179**, 513–533 (2016).
- 518 17. Jamil, M. *et al.* Using multiple machine learning algorithms to classify elite and sub-
519 elite goalkeepers in professional men’s football. *Sci. Rep.* **11**, 22703 (2021).
- 520 18. Barron, D., Ball, G., Robins, M. & Sunderland, C. Artificial neural networks and
521 player recruitment in professional soccer. *PLoS One* **13**, e0205818 (2018).
- 522 19. García-Aliaga, A., Marquina, M., Coterón, J., Rodríguez-González, A. & Luengo-
523 Sánchez, S. In-game behaviour analysis of football players using machine learning
524 techniques based on player statistics. *Int. J. Sports Sci. Coach.* **16**, 148–157 (2021).
- 525 20. Oliver, J. L. *et al.* Using machine learning to improve our understanding of injury risk
526 and prediction in elite male youth football players. *J. Sci. Med. Sport* **23**, 1044–1048
527 (2020).
- 528 21. Rommers, N. *et al.* A machine learning approach to assess injury risk in elite youth
529 football players. *Med. Sci. Sports Exerc.* **52**, 1745–1751 (2020).
- 530 22. Ruddy, J. D. *et al.* Modeling the risk of team sport injuries: a narrative review of
531 different statistical approaches. *Front. Physiol.* **10**, (2019).
- 532 23. McCulloch, W. S. & Pitts, W. A logical calculus of the ideas immanent in nervous
533 activity. *Bull. Math. Biophys.* **5**, 115–133 (1943).
- 534 24. Akyildiz, Z. *et al.* Monitoring the post-match neuromuscular fatigue of young Turkish

- 535 football players. *Sci. Rep.* **12**, 13835 (2022).
- 536 25. Houtmeyers, K. C., Jaspers, A. & Figueiredo, P. Managing the training process in elite
537 sports: from descriptive to prescriptive data analytics. *Int. J. Sports Physiol. Perform.*
538 **16**, 1719–1723 (2021).
- 539 26. Bzdok, D., Altman, N. & Krzywinski, M. Points of Significance: Statistics versus
540 machine learning. *Nat. Methods* **15**, 233–234 (2018).
- 541 27. Richter, C., O'Reilly, M. & Delahunt, E. Machine learning in sports science:
542 challenges and opportunities. *Sport. Biomech.* 1–7 (2021)
543 doi:10.1080/14763141.2021.1910334.
- 544 28. Tian, T., Song, C., Ting, J. & Huang, H. A French-to-English machine translation
545 model using transformer network. *Procedia Comput. Sci.* **199**, 1438–1443 (2022).
- 546 29. Bhandari, H. N. *et al.* Predicting stock market index using LSTM. *Mach. Learn. with*
547 *Appl.* **9**, 100320 (2022).
- 548 30. Ettensperger, F. Comparing supervised learning algorithms and artificial neural
549 networks for conflict prediction: performance and applicability of deep learning in the
550 field. *Qual. Quant.* **54**, 567–601 (2020).
- 551 31. Anaby-Tavor, A. *et al.* Do not have enough data? Deep learning to the rescue! *Proc.*
552 *AAAI Conf. Artif. Intell.* **34**, 7383–7390 (2020).
- 553 32. Rossi, A. *et al.* Effective injury forecasting in soccer with GPS training data and
554 machine learning. *PLoS One* **13**, e0201264 (2018).
- 555 33. Goddard, J. Regression models for forecasting goals and match results in association
556 football. *Int. J. Forecast.* **21**, 331–340 (2005).
- 557 34. Kass, G. V. An exploratory technique for investigating large quantities of categorical
558 data. *Appl. Stat.* **29**, 119 (1980).
- 559 35. Hunt, E. B., Marin, J. & Stone, P. J. . *Experiments in induction.* Academic Press
560 (1966).
- 561 36. Gomes, H. M. *et al.* Adaptive random forests for evolving data stream classification.
562 *Mach. Learn.* **106**, 1469–1495 (2017).
- 563 37. Cortes, C. & Vapnik, V. Support-vector networks. *Mach. Learn.* **20**, 273–297 (1995).
- 564 38. Murty, M. N., Raghava, R., Murty, M. N. & Raghava, R. Kernel-based SVM. Support
565 vector machines and perceptrons: Learning, optimization, classification, and
566 application to social networks. in 57–67 (2016).
- 567 39. Bishop, C. M. *Neural networks for pattern recognition.* (Oxford University
568 PressOxford, 1995). doi:10.1093/oso/9780198538493.001.0001.

- 569 40. Ruddy, J. *et al.* Predictive modeling of hamstring strain injuries in elite Australian
570 footballers. *Med. Sci. Sport. Exerc.* **50**, 906–914 (2018).
- 571 41. Chen, Y. *et al.* Auto-encoding variational Bayes. *Cambridge Explor. Arts Sci.* **2**,
572 (2024).
- 573 42. Goodfellow, I. *et al.* Generative adversarial networks. *Commun. ACM* **63**, 139–144
574 (2020).
- 575 43. Karim, M. R. *et al.* Deep learning-based clustering approaches for bioinformatics.
576 *Brief. Bioinform.* **22**, 393–415 (2021).
- 577 44. Xu, Q. & He, X. Football training evaluation using machine learning and decision
578 support system. *Soft Comput.* **26**, 10939–10946 (2022).
- 579 45. Mnih, V. *et al.* Human-level control through deep reinforcement learning. *Nature* **518**,
580 529–533 (2015).
- 581 46. Silver, D. *et al.* Mastering the game of Go with deep neural networks and tree search.
582 *Nature* **529**, 484–489 (2016).
- 583 47. Wheelen, T., Hoffman, J. & Bamford, C. Strategic management and business policy
584 globalization. in *Pearson* (2018).
- 585 48. Nassis, G. P. Leadership in science and medicine: can you see the gap? *Sci. Med.*
586 *Footb.* **1**, 195–196 (2017).
- 587 49. Gabbett, H. T., Windt, J. & Gabbett, T. J. Cost-benefit analysis underlies training
588 decisions in elite sport. *Br. J. Sports Med.* **50**, 1291–1292 (2016).
- 589 50. Edouard, P., Verhagen, E. & Navarro, L. Machine learning analyses can be of interest
590 to estimate the risk of injury in sports injury and rehabilitation. *Ann. Phys. Rehabil.*
591 *Med.* **65**, 101431 (2022).

592

593

594 **Authors contribution**

595 All authors contributed to the writing of the paper. All authors read and approved the final
596 version.

597

598 **Conflict of interest**

599 The authors declare no conflict of interest for this paper.

600

601 **Data availability statement**

602 This manuscript does not have associated data.

603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618

Table 1. Supervised and unsupervised machine learning analysis.

<p>Regression</p>	<p>Supervised ML regression is a type of predictive analysis that is used to model and analyze relationships between variables. It aims to predict a continuous target variable based on one or more independent variables. The goal is to find the best fit line or curve that minimizes the difference between predicted and actual values. This is achieved through algorithms that adjust the weights of input features to reduce error in predictions. Regression techniques are widely used in fields such as finance, medicine, and environmental science for tasks like prediction market value and estimating injury risk.</p>	<p>Examples of regressions analysis: <i>Boosting, Decision Tree, K Nearest Neighbor, Neural Network, Random Forest, Regularized linear, Support Vector Regression</i></p>
<p>Classification</p>	<p>Supervised ML classification is a type of algorithm used to assign predefined labels to new data points. It works by learning from a dataset with known labels and then applying this knowledge to categorize new, unlabeled data. Common applications include sport movements analysis and medical diagnosis, where the algorithm must decide which category or class the new data belongs to based on its features.</p>	<p>Examples of classifications analysis: <i>Boosting, Decision Tree, K Nearest Neighbor, linear discriminant, Neural Network (includes Deep CNNs), Random Forest, Support Vector Machine.</i></p>
<p>Clustering</p>	<p>Clustering in ML is an unsupervised learning technique used to group a set of objects in such a way that objects in the same group (called a cluster) are more similar to each other than to those in other groups (clusters). It is commonly used in statistical data analysis for pattern recognition, game-tactical analysis, information retrieval, and bioinformatics. Algorithms like K-Means, Hierarchical clustering, and DBSCAN are popular methods for performing clustering tasks. The goal is to discover the inherent structure within the data, often to identify distinct subgroups without pre-labeled data or human supervision.</p>	<p>Examples of clustering analysis: <i>Density-based, Fuzzy C-Means, Hierarchical, Neighborhood-based.</i></p>

Table 2. This table reports the general characteristics of a SWOT analysis. The SWOT analysis process serves as a compass, guiding organizations toward effective strategies, risk management, and sustainable growth of a business (club).

Strategic planning	SWOT analysis helps organizations develop effective strategies by identifying their strengths, weaknesses, opportunities, and threats. It provides a comprehensive view of the current situation, enabling informed decision-making.
Self-reflection and awareness	Organizations need to understand their internal capabilities (strengths and weaknesses) and external factors (opportunities and threats). SWOT analysis encourages self-reflection and awareness, leading to better alignment with organizational goals.
Risk assessment	By evaluating potential threats (such as market changes, competition, or regulatory issues), organizations can proactively address risks. SWOT analysis allows them to prioritize risk mitigation strategies.
Resource allocation	SWOT analysis guides resource allocation. Organizations can allocate resources more effectively by capitalizing on strengths and minimizing weaknesses. It helps prioritize investments and efforts
Competitive advantage	Identifying unique strengths and opportunities allows organizations to create a competitive edge. Leveraging these advantages helps them stand out in the market.
Adaptation to change	The business landscape constantly evolves. SWOT analysis enables organizations to adapt to changes by recognizing emerging opportunities and addressing potential threats promptly
Communication and alignment	SWOT analysis fosters communication among team members, stakeholders, and leadership. It aligns everyone around a common understanding of the organization's position and future direction.

Figure 1: An example of SWOT analysis regarding the use of ML for injury risk assessment for a football team