

# New Machine Learning Approach for Detection of Injury Risk Factors in Young Team Sport Athletes

## Authors

Susanne Jauhiainen<sup>1</sup>✉, Jukka-Pekka Kauppi<sup>1</sup>, Mari Leppänen<sup>2</sup>, Kati Pasanen<sup>2, 3, 4, 5</sup>, Jari Parkkari<sup>2, 6</sup>, Tommi Vasankari<sup>2</sup>, Pekka Kannus<sup>2, 6</sup>, Sami Äyrämö<sup>1</sup>

## Affiliations

- 1 Faculty of Information Technology, University of Jyväskylä, Jyväskylä, Finland
- 2 Tampere Research Centre of Sports Medicine, UKK Institute, Tampere, Finland
- 3 Sport Injury Prevention Research Centre, Faculty of Kinesiology, University of Calgary, Calgary, Alberta, Canada
- 4 Alberta Children's Hospital Research Institute, University of Calgary, Calgary, Alberta, Canada
- 5 McCaig Institute for Bone and Joint Health, University of Calgary, Calgary, Alberta, Canada
- 6 Tampere University Hospital, Tampere, Finland

## Key words

sports medicine, predictive methods, machine learning, knee injuries, ankle injuries, basketball and floorball

accepted 20.07.2020

Published online: 2020

## Bibliography

Int J Sports Med

DOI 10.1055/a-1231-5304

ISSN 0172-4622

© 2020. Thieme. All rights reserved.

Georg Thieme Verlag KG, Rüdigerstraße 14,  
70469 Stuttgart, Germany

## Correspondence

Susanne Jauhiainen

Faculty of Information Technology, University of Jyväskylä

Mattilanniemi 2

40100 Jyväskylä

Finland

Tel.: 358504043416

susanne.m.jauhiainen@jyu.fi

## ABSTRACT

The purpose of this article is to present how predictive machine learning methods can be utilized for detecting sport injury risk factors in a data-driven manner. The approach can be used for finding new hypotheses for risk factors and confirming the predictive power of previously recognized ones. We used three-dimensional motion analysis and physical data from 314 young basketball and floorball players (48.4 % males,  $15.72 \pm 1.79$  yr,  $173.34 \pm 9.14$  cm,  $64.65 \pm 10.4$  kg). Both linear (L1-regularized logistic regression) and non-linear methods (random forest) were used to predict moderate and severe knee and ankle injuries ( $N = 57$ ) during three-year follow-up. Results were confirmed with permutation tests and predictive risk factors detected with Wilcoxon signed-rank-test ( $p < 0.01$ ). Random forest suggested twelve consistent injury predictors and logistic regression twenty. Ten of these were suggested in both models; sex, body mass index, hamstring flexibility, knee joint laxity, medial knee displacement, height, ankle plantar flexion at initial contact, leg press one-repetition max, and knee valgus at initial contact. Cross-validated areas under receiver operating characteristic curve were 0.65 (logistic regression) and 0.63 (random forest). The results highlight the difficulty of predicting future injuries, but also show that even with models having relatively low predictive power, certain predictive injury risk factors can be consistently detected.

## Introduction

Sports injuries are very common across different sports among both elite and recreational athletes [1–3]. They can have significant effects on health and performance and may even cause prolonged problems in a person's life [3]. Sports injuries can lead to pain, loss of playing or working time, and decreased motility and stability [3].

The incidence rate of some injuries, such as the anterior cruciate ligament (ACL) injury, is a growing cause for concern [4]. Effective prevention of injuries presumes that the most relevant risk factors are found. Even though many intrinsic and extrinsic risk factors have been identified, there is no clear consensus on the findings [5].

A large majority of existing sports injury studies rely on an explanatory analysis approach [6, 7]. Explanatory methods have played an important role in the development of sports injury research and will be needed in future research as well. They are used when the purpose is to explain or understand data or phenomena of interest. However, high explanatory power does not necessarily imply high predictive power [8]. Therefore, risk factors that are identified by explanatory methods demonstrate only a statistically significant association with injuries but may not have predictive power on them [6, 7].

Another limitation of explanatory analyses is that they often focus on a small number of variables and their linear associations with injuries in isolation. However, underlying causes behind sports injuries have been considered to be multifactorial, indicating that a high number of variables and their inter-relationships should be considered [9, 10]. It has also been suggested that using cut-off values and studying only linear interactions between isolated variables cannot successfully identify injury predictors, but more complex models should be applied [11]. To overcome these limitations, predictive analysis should be utilized alongside explanatory methods. This has been previously suggested specifically for sports injury research as well [12].

Predictive analysis focuses on predicting new or future observations from data [8]. By exploiting computational power, predictive methods are able to analyze a larger set of variables including their interactions and nonlinear relationships as well as to efficiently remove redundant variables from a model. Therefore, they can be used for generating new hypotheses for sports injury risk factors in a data-driven manner.

In predictive analysis, the generalization ability of a model should always be assessed on independent test data, i. e., data that have not been used in the training phase. This measures how accurate the trained model will be on new unseen observations, and only after such validation can any conclusions about the predictive power be drawn [8]. In addition, when constructing a predictive model, it is necessary to confirm that the prediction results are significantly above the random chance level. This kind of confirmatory analysis is especially relevant with smaller sample sizes. If this issue is not considered, in the worst case it can lead to false interpretations and conclusions. In neuroscience, for example, the problem has been widely recognized [13]. Permutation tests can be used to confirm the significance of the models and relevance of the chosen predictors [13].

Another important issue related to predictive analysis is the explainability of a model. Explainability means that the model somehow explains its predictions, for example, gives information on how individual variables contribute to the prediction outcome rather than predicting as a black box [14]. Explainable models and their predictions are more informative, easier to trust, and therefore can provide more practical benefits. A term widely used with sophisticated machine learning methods is explainable artificial intelligence (XAI) [14]. In some domains, such as medicine, model explainability is considered highly important [15] and should be pursued in sports science and medicine as well.

Over the last few years, the first studies using predictive analysis in sports injury research were conducted [6, 9, 16, 17]. Previous studies have, however, focused solely on the prediction task with-

out paying attention to the explainability of the models. In addition, two of the studies also used a very low number of variables (from 3–11), although a larger set might have increased the accuracy [9, 16]. The need and potential of predictive machine learning methods in sports injury prediction have been recognized, but more research is needed [12, 17].

Therefore, the aim of this study is to utilize predictive machine learning methods to detect variables with predictive power on sports injuries. We present a framework that can be used to detect consistent injury predictors in a data-driven manner and validate their predictive power on independent test data. Consistent means that the variable is constantly chosen as an important predictor in the model used. Our framework utilizes both linear and non-linear classification methods, namely L1-regularized logistic regression and random forests, to predict moderate and severe knee and ankle injuries. Generalization ability of these models is assessed with 10-fold cross-validation. A reference model based on randomized labels is constructed to confirm that the observed prediction performance is not achieved by chance. Consistent injury predictors are detected with a Wilcoxon signed-rank test. This approach can be used for finding new hypotheses for injury risk factors as well as confirming the predictive power of previously recognized risk factors. Our secondary aim is to compare linear and non-linear methods for the task.

## Materials and Methods

### Participants

The data were collected in the Predictors of Lower Extremity Injuries in Team Sports (PROFITS) study [18]. The study was conducted in accordance with the Declaration of Helsinki and was approved by the Ethics Committee of the Pirkanmaa Hospital District, Tampere, Finland (ETL-code R10169). The authors declare that this study meets the ethical requirements of the journal [19]. Altogether 175 basketball and 139 floorball youth (12–21 years) players, including 162 females ( $15.44 \pm 1.95$  yrs,  $167.92 \pm 6.44$  cm,  $60.86 \pm 8.58$  kg) and 152 males ( $16.03 \pm 1.59$  yrs,  $179.13 \pm 8.00$  cm,  $68.68 \pm 10.76$  kg) from the two highest junior league levels of the Tampere city district, Finland, were recruited. To be included they had to be official team members (i. e., have valid playing contract and licenses), 21 years old or younger at baseline, and free from injury at baseline. Information about previous injuries, their treatment, and whether the player was fully recovered were assessed with a baseline questionnaire. The players entered the study during the preseason in 2011, 2012, or 2013. They signed a written informed consent form before inclusion (including parental consent for players aged  $\leq 18$  years).

### Data collection

At baseline, each player participated in physical tests including a vertical drop jump (VDJ) (3D motion analysis), height, weight, isokinetic concentric quadriceps and hamstring strength, isometric hip abductor strength, one repetition maximum (1RM) leg press, knee joint laxity (KT-1000), generalized joint laxity (Beighton scale), genu recurvatum, navicular drop, hip anteversion, and hamstring

flexibility (for more details, see Supplementary ► **Table 1S** and on-line supplementary appendices in [18]).

The VDJ was performed from a 30-cm box. Players were instructed to drop off the box and perform a maximal jump upon landing with their feet on two separate force platforms (BP6001200; AMTI, Watertown, MA, USA). The 3D motion analysis was carried out using sixteen reflective markers placed over anatomic landmarks on the lower extremities according to the Plug-In Gait Marker set (Vicon Nexus v. 1.7; Oxford Metrics, Oxford, UK) and eight high-speed cameras (Vicon T40). Kinetics and kinematics variables were extracted using the Vicon Nexus Plug-in Gait model. Medial knee displacements were extracted using a custom MATLAB script (MathWorks, Inc., Natick, MA, USA). For more detailed description of the motion data collection and variable extraction, see [18, 20].

The injury definition was based on the time-loss definition by Fuller et al. [21]. We focused on moderate to severe acute non-contact knee and ankle injuries that resulted in an athlete being unable to fully participate in training or match play for at least 8 days. Non-contact injury was defined as an injury that occurred without direct contact to the injured body part. Injuries were recorded by a team coach or another designated team member. For injury registration, the study physicians contacted the team coach or designated team member on a weekly basis by phone or email. The designated team member was someone who was always present at practice and matches, e.g., head, assistant, or strength and conditioning coach, team manager, or physiotherapist. The study physicians contacted the athlete after each injury and collected information about the injury time, place, cause, type, location, and the time-loss due to the injury in a standardized phone interview. For exposure registration, the team coaches recorded player participation in team practice and game play and emailed the records to the study group at the end of each month.

## Data preprocessing

All data analysis was performed with MATLAB R2016b (MathWorks, Inc) and classification methods run with the Statistics and Machine Learning Toolbox 11.0. For classification, the players with moderate and severe acute ankle and knee injuries formed the first group (group A,  $n = 57$ ) and players with no injuries formed the other (group B,  $n = 257$ ). Athletes with mild injuries (time-loss  $\leq 7$  days,  $n = 21$ ) were excluded from the analysis. Altogether 58 variables were chosen for further analysis by a group of experts in sports medicine, including a sports medicine researcher and four clinical researchers (one physiotherapist and three physicians). Four variables had more than 50 % of missing values (iliopsoas and quadriceps extensibility from both legs) because they were added to the test patterns only in the second year of testing. They were excluded from the analysis, resulting into 54 variables. The chosen variables are described in the Supplementary ► **Table 1S**.

After exclusion of the irrelevant and sparse variables, 22 variables with missing data remained and were imputed with K-nearest neighbor imputation with a  $k$  value of 10. On average, each of these 22 variables had five missing values (1.6 % of the 314 observations). Data was normalized to have a mean of zero and standard deviation of one for each column. The variables that had been measured separately for both right and left legs were transformed to dominant (leg used for kicking a ball) and non-dominant leg variables.

## Choice of classification methods

Two commonly used methods, random forest and L1-regularized logistic regression, were chosen for the binary classification task in our framework. These methods were selected because of their in-built variable importance features. Random forest is a nonlinear classification and regression method that has become a standard data analysis tool in different fields such as medicine and bioinformatics [22] and has been used in sports injury research as well [23]. It is based on building an ensemble of multiple decision trees [24]. The model was trained with a hundred trees [24] and the minimum number of observations per tree leaf and the number of predictors to sample at each split were chosen using Bayesian optimization. To estimate the predictive power of the variables, we recorded and analyzed the out-of-bag estimates of variable importance [24].

L1-regularized logistic regression, in turn, is a linear classification method that has been used to model sports injury outcomes [23]. One benefit of this method is that it is capable of automatically discarding redundant and/or irrelevant variables from the model. This is done by penalizing the model with the L1 norm and as a result, some of the variable coefficients tend to shrink to exactly zero. The optimal amount of penalization was estimated with stratified 10-fold cross-validation.

Variable importance for logistic regression was based on the variable coefficient values. We analyzed whether a variable was chosen as a predictor in the model, i.e., the variable coefficient was not shrunk to zero. Variable importance was then the number of times the variable was chosen over the ten CV folds (a value between zero and ten). The sign of each variable coefficient was also assessed in order to perceive whether the variable decreased or increased the injury risk.

## Validation

Generalization ability of our models was assessed with 10-fold cross-validation (CV). K-fold CV is based on randomly splitting the data into K sets and leaving each set at a time for testing while the rest of the sets are used to train a model. Test performance was assessed with Area Under the Receiver Operating Characteristics Curve (AUC-ROC) [25]. It is based on both true positive and false positive rates and it can be used with imbalanced class distributions, which is the case in our data. AUC-ROC provides a value of 1.0 for perfect prediction and 0.5 for purely random prediction.

AUC-ROC and variable importance values were estimated by ten-fold cross-validation. Normalization and imputation of the training data were done separately inside each fold and the test data were then normalized using coefficients estimated from the training data. Because K-fold CV is based on random splitting of the data, there is variation in the K-fold validation estimates [26]. Therefore, the analysis was repeated a hundred times and results were averaged over the runs to obtain a more reliable estimate for the generalization ability.

## Confirmatory data analysis

To confirm the significance of our results, permutation tests were used [13]. A reference model was constructed by randomly shuffling the class labels in the training data. By comparing the outcome of the true models to the distribution of values from the random models, we confirmed that the performance was not observed by

chance. In addition, we can detect significantly consistent injury predictors by comparing the variable importance of the true and the random reference models. If a variable is consistently important in the true model but not in the reference model, that confirms its significance in the prediction.

To confirm the significance of obtained performance, a paired comparison between AUC-ROC values of the true and random model from a hundred repeated 10-fold CV runs was conducted based on a Wilcoxon signed-rank test. In each CV run, the fold divisions were kept the same for random and true models to allow fair pairwise comparison.

To detect significantly consistent injury predictors, we compared the variable importance values. Again, the values from the hundred repetitions were compared between the random and true models but with a Wilcoxon signed-rank test. The limit of significance was set to  $\alpha = 0.01$  and corrected with Bonferroni correction. The framework used is summarized in ► Fig. 1.

## Results

### Random forest

Random forest suggested twelve consistent injury predictors ( $p < 0.01$ ). The variable importance values averaged over the CV folds and a hundred repeated runs can be seen in ► Fig. 2. The larger the importance value, the greater the importance of the variable is for the prediction task. By comparing the values between true and randomized results, variables with true predictive power can be detected. If the value of the true model is significantly larger

than the value of random model, its predictive power is not likely the result of chance or noise in data. Negative values indicate the variable was not important in the prediction.

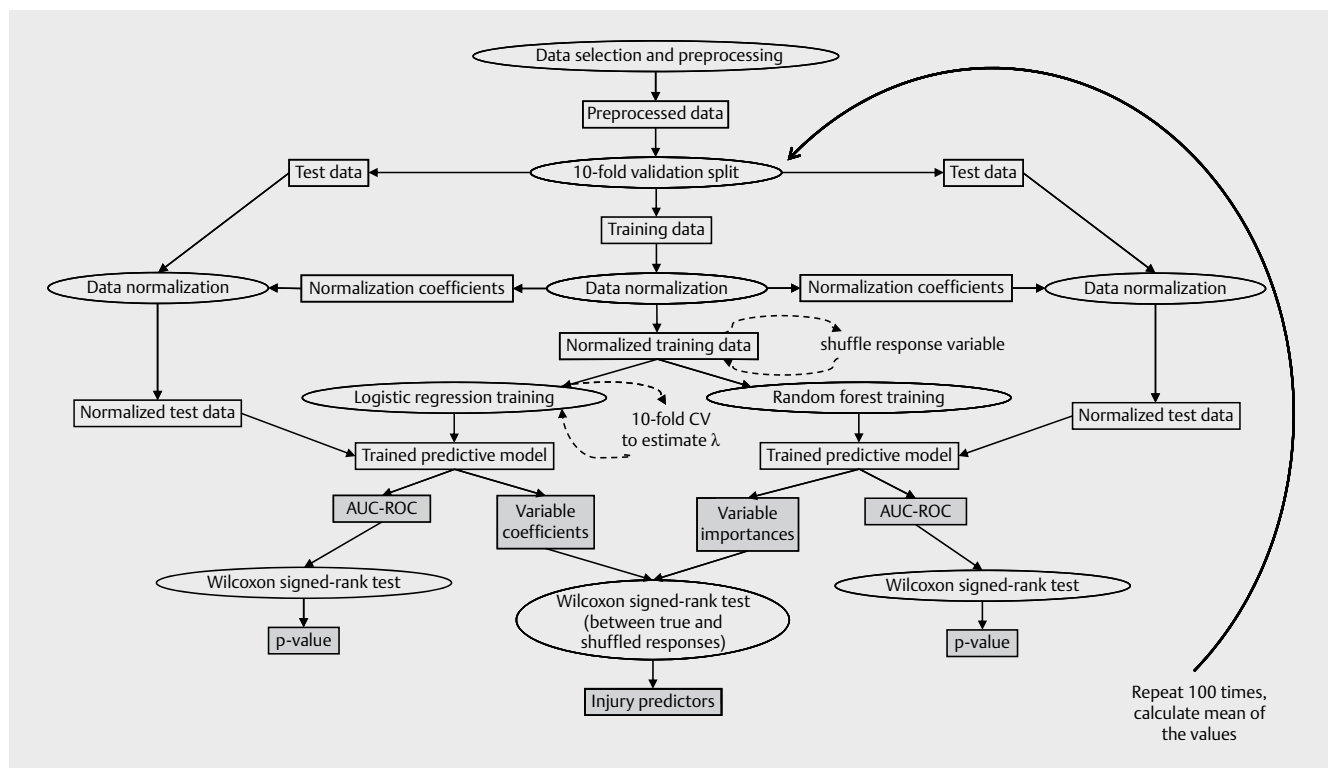
As seen in ► Fig. 2, sex, hamstring flexibility (both dominant and non-dominant legs), body mass index (BMI), KT1000 (dominant leg), and height show the highest random forest importance values. Other suggested predictors include leg press 1RM, knee valgus at IC (dominant leg), knee flexion peak (non-dominant leg), medial knee displacement (dominant leg), ankle flexion at IC (dominant leg), and navicular drop (non-dominant leg).

The mean AUC-ROC value for random forest was 0.63 (0.94 for the training data). The AUC-ROC values were higher ( $p < 0.001$ ) with real responses than the randomized ones (mean AUC-ROC 0.48), which confirms the significance of the random forest models.

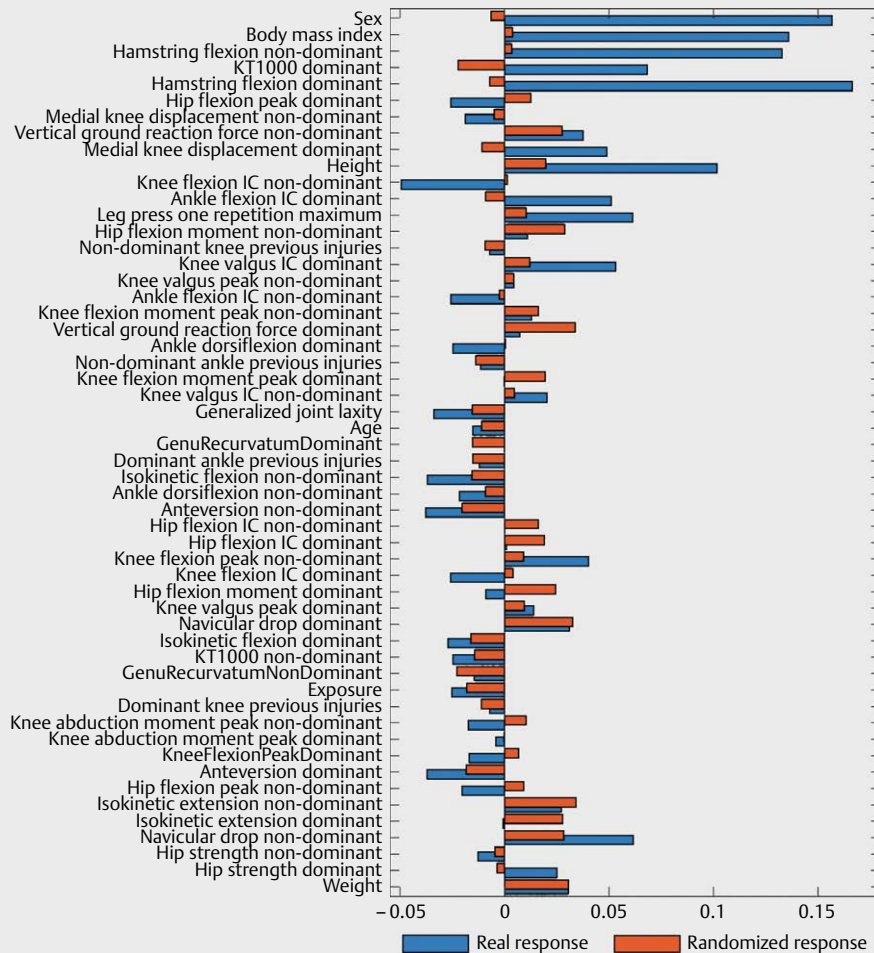
### Logistic regression

► Figure 3 shows the variables chosen most frequently as predictors in the L1-regularized logistic regression. The bars represent the number of CV folds where a variable was chosen for the predictive model (i. e., its coefficient was not shrunk to zero). As can be seen in the figure, some of the variables were chosen for prediction in almost every CV split, whereas the others were regarded as not important and their coefficients shrunk to exactly zero. Twenty variables were suggested as consistent injury predictors ( $p < 0.01$ ) with the logistic regression model.

The suggested variables were sex, BMI, hamstring flexibility (both legs), KT1000 (dominant leg), hip flexion peak (dominant leg), medial knee displacement (both legs), vertical ground reaction force (vGRF) (both legs), height, knee flexion at IC (non-dom-



► Fig. 1 Framework of the proposed predictive analysis approach.



► **Fig. 2** Variable importance values from random forest. Blue bars correspond to the results with real response, red ones with randomized response.

inant leg), ankle flexion at IC (both legs), leg press 1RM, hip flexion moment peak (non-dominant leg), previous injuries of non-dominant knee, knee valgus at IC (dominant leg), knee valgus peak (non-dominant leg), and knee flexion moment peak (non-dominant leg). In the figure, these are the twenty variables with the highest frequency value.

The mean AUC-ROC value for logistic regression models was 0.65 (0.76 for the training data). The AUC-ROC values were higher ( $p < 0.001$ ) with real responses than the randomized ones (mean AUC-ROC 0.50).

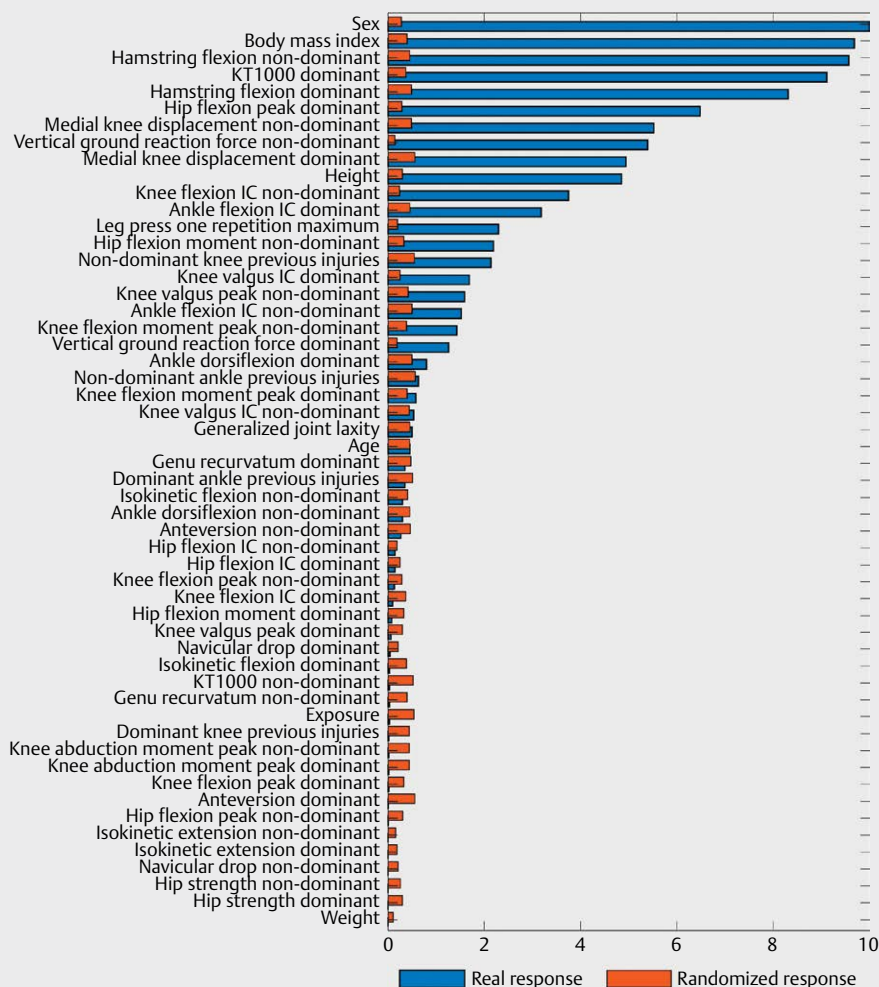
### Logistic regression coefficients

Whenever a variable was chosen for the logistic regression model, the direction of the coefficient was extremely consistent, always either positive or negative. Therefore, over all the folds and a hundred runs, the variable always had a similar effect on the prediction, i. e., it either increased or decreased the risk of injury. Directions of variable coefficients for the ten most often selected variables as well as those that were found by both models can be seen in ► **Table 1**.

Based on the coefficients, female sex contributes to a greater risk than male (male = 1, female = 2 in data) as well as larger BMI, lower height, and higher leg press 1RM result. Higher hamstring flexibility and vGRF of both legs increase the risk of injury. The higher value of KT1000 of the dominant leg as well as higher hip flexion peak and knee flexion at IC of the non-dominant leg also contribute to the injury risk. Less ankle plantar flexion (negative values) and larger knee valgus angle (negative values) of the dominant leg contribute to the higher risk. Interestingly, for medial knee displacement, the direction was different between the legs. For the non-dominant leg, higher medial knee displacement increased the risk, but for the dominant leg, a lower value increased it.

### Consistent injury predictors chosen by both methods

The following ten variables were suggested as consistent injury predictors ( $p < 0.01$ ) by both models: sex, body mass index, hamstring flexibility (non-dominant leg), KT1000 (dominant leg), hamstring flexibility (dominant leg), medial knee displacement (dominant leg), height, ankle (plantar) flexion at IC (dominant leg), leg press one repetition maximum (1RM), and knee valgus at IC (dominant leg).



► **Fig. 3** Variable importances for L1-regularized logistic regression. Measured as the number of times each variable was chosen over the ten CV folds. Blue Dark bars correspond to the results with real response, red to the randomized response

## Discussion

The purpose of this study was to utilize predictive machine learning methods to detect variables with predictive power for sports injuries. Multiple injury risk factors have been recognized in previous explanatory studies, but the predictive power of these variables remains unclear until tested on independent data. We presented a framework that detects consistent injury predictors in a data-driven manner and validates their predictive power on independent test data. This approach can be used for finding new hypotheses for injury risk factors as well as confirming the predictive power of previously recognized risk factors. Any new hypotheses should then be confirmed by domain experts in future studies, including utilizing explanatory methods.

Despite the low predictive accuracy (AUC = 0.65), a set of ten consistent injury predictor variables was detected by both models. The obtained AUC score is in line with the previous studies [6, 9, 16, 17] and confirms the difficulty of predicting sports injuries. A recently published predictive analysis study that compared different methods and their injury prediction accuracies obtained

an AUC score of 0.747 when predicting lower extremity muscle injuries in 132 male professional soccer and handball players [9]. A paper by Dower and colleagues [17] utilized time series data and artificial neural networks, achieving AUC scores between 0.75 and 0.80 on average when predicting soft tissue injuries in Australian football players.

Another study found that previously detected risk factors with explanatory power had a very poor predictive performance (median AUC scores 0.57 and 0.52) on hamstring strain injuries in 362 elite Australian footballers [16]. However, this study used a small number of variables in the prediction (three and eight). In addition, previous studies have focused solely on the prediction task without considering the explainability of the predictive model. Explainable models that assess the effect of each variable in prediction, for example, are easier to trust and provide more practical information to the domain experts.

Most of the injury predictor variables suggested in our study are supported by previous research. Our results suggest that female sex, larger BMI, and lower height increased the risk of acute non-contact knee and ankle injury. Previous explanatory research has



► **Table 1** The number of coefficients with positive, negative, and zero values over the ten folds and one hundred runs.

Variable	Positive	Negative	Zero
Sex	0	999	1
Body mass index	968	0	32
Hamstring flexion non-dominant	957	0	43
KT1000 dominant	911	0	89
Hamstring flexion dominant	831	0	169
Hip flexion peak dominant	648	0	352
Medial knee displacement non-dominant	552	0	448
Vertical ground reaction force non-dominant	539	0	461
Medial knee displacement dominant	0	494	506
Height	0	485	515
Knee flexion IC non-dominant	375	0	625
Ankle flexion IC dominant	318	0	682
Leg press one repetition maximum	230	0	770
Knee valgus IC dominant	0	169	831
Vertical ground reaction force dominant	126	0	874

detected similar associations with lower extremity sports injuries [2, 5, 27, 28]. For muscle flexibility, there are contradictory findings [5, 29]. Our results suggest that increased hamstring flexibility of both the dominant and non-dominant leg contributes to a higher risk of acute non-contact knee and ankle injuries.

Concerning the association between muscle strength and sports injury risk, the findings are conflicting [30, 31]. Our study found higher leg press 1RM to be associated with higher injury risk. This could be, for example, because stronger athletes exert greater forces and moments to the joints and muscles during activity, are more mature, and tend to train more and perform at higher levels. Also, our findings that increased knee laxity (KT-1000) and less ankle plantar flexion at IC of the dominant leg contribute to higher injury risk have been previously recognized [32, 33].

Our results suggest that larger knee valgus and medial knee displacement of the non-dominant leg increase the risk of acute non-contact knee and ankle injury. Associations between knee valgus loading and risk of lower extremity injuries have been found previously [34]. However, our results also suggested that smaller medial knee displacement of the dominant leg increased the risk, which contradicts the results of the non-dominant leg. In the group of non-injured athletes, the medial knee displacement of the dominant leg is notably larger than in the non-dominant leg. In the injured group, there is no such difference (see Supplementary ► **Table 1S**). This difference between sides is causing the conflicting regression coefficients inside the framework. However, such differences were not observed in the knee valgus angles. This may be because the medial knee displacement is more sensitive to the athlete rotating during landing. In our data, approximately 74 % of the athletes rotated towards the side of their dominant leg during VDJ. Another possible explanation may simply be differences in the use of the dominant and non-dominant leg.

Our secondary aim was to assess differences between linear and non-linear methods. In our prediction task, the predictive accuracy of the linear L1-regularized logistic regression was slightly better (AUC = 0.65) than the accuracy of the non-linear random forest model (AUC = 0.63). The difference, however, is negligible for drawing conclusions of their mutual superiority. The suggested injury risk factors were largely the same for both models, but logistic regression suggested a larger set of predictors. Generally, we believe it can be beneficial to utilize a combination of methods to detect the most relevant injury risk factors.

The strength of our approach is that with predictive methods and confirmatory analysis, consistent injury predictors can be detected even from data with weak phenomena. For example, with small datasets the approach can help to avoid findings based on chance. Thus, it can be useful in other sports science and medicine studies as well, even though the used data itself does not necessarily possess high predictive power or strong phenomena. Another strength is the prospective data collection of a large number of variables from a large cohort of athletes. Predictive methods utilize computational power and thus enable analysis of all relevant data and do not require exclusion based on prior assumptions. In addition, our study uses a well-defined prediction outcome of moderate and severe knee and ankle injuries whose risk factors have been previously established in explanatory research.

However, there are also limitations related to the data used. After baseline data was collected, the injury follow-up lasted for 12 months. Many of the collected variables may, however, change notably during this period, especially in young athletes [10]. In the future, more comprehensive data that observes short-term changes in variables should be collected because there can be changes, for example, based on the time in season and weekly training and game loads. Wearable technologies, for example, allow continuous monitoring of athletes. It can be expected that time series data from wearable devices combined with applicable predictive methods will increase the prediction accuracy as the study by Dower et al. indicated [17].

To conclude, in order to have practical value in the clinical assessment of injury risk, the predictive accuracy of the presented models that were trained on the prospective data should be improved. The models were, however, able to detect a set of consistent injury predictors. Thus, the approach can be useful for finding new hypotheses for injury risk factors as well as confirming the predictive power of risk factors found in previous explanatory studies. Although the achieved predictive accuracy of our study remained relatively low (AUC = 0.65), a set of ten consistent injury predictor variables was detected by both models (sex, body mass index, hamstring flexibility, knee joint laxity, medial knee displacement, height, ankle plantar flexion at initial contact, leg press one-repetition max, and knee valgus at initial contact). The obtained accuracy is in line with previous studies and confirms that predicting sports injuries is a cumbersome task. More research is required to find risk factors that best predict injury and should include more comprehensive data. The obtained performance was similar between the linear and non-linear methods. Future research is needed to assess the suitability and performance of linear versus non-linear methods in sports injury prediction tasks.

## Funding

This study was supported by the Finnish Ministry of Education and Culture, and Competitive State Research Financing of the Expert Responsibility area of Tampere University Hospital (grants 9S047, 9T046, 9U044, 9N053). This work has been carried out in two projects "Value from health data with cognitive computing" and "Watson Health Cloud", funded by Business Finland. Susanne Jauhiainen was funded by the Jenny and Antti Wihuri Foundation (grant 00180121). Jukka-Pekka Kauppi was funded by the Academy of Finland Postdoctoral Researcher program (Research Council for Natural Sciences and Engineering; grant 286019).

## Conflict of Interest

The authors declare that they have no conflict of interest.

## References

- Jacobsson J, Timpka T, Kowalski J et al. Prevalence of musculoskeletal injuries in Swedish elite track and field athletes. *Am J Sports Med* 2012; 40: 163–169
- Emery CA, Rose MS, McAllister JR et al. A prevention strategy to reduce the incidence of injury in high school basketball: a cluster randomized controlled trial. *Clin J Sport Med* 2007; 17: 17–24
- Myklebust G, Holm I, Mæhlum S et al. Clinical, functional, and radiologic outcome in team handball players 6–11 years after anterior cruciate ligament injury: A follow-up study. *Am J Sports Med* 2003; 31: 981–989
- Bahr R, Holme I. Risk factors for sports injuries – a methodological approach. *Br J Sports Med* 2003; 37: 384–392
- Murphy DF, Connolly DAJ, Beynnon BD. Risk factors for lower extremity injury: A review of the literature. *Br J Sports Med* 2003; 37: 13–29
- Rossi A, Pappalardo L, Cintia P et al. Effective injury forecasting in soccer with GPS training data and machine learning. *PLoS One* 2018; 13: e0201264
- Bahr R. Why screening tests to predict injury do not work – and probably never will....: A critical review. *Br J Sport Med* 2016; 50: 776–780
- Shmueli G. To explain or to predict? *Stat Sci* 2010; 25: 289–310
- López-Valenciano A, Ayala F, Puerta JM et al. A preventive model for muscle injuries: A novel approach based on learning algorithms. *Med Sci Sports Exerc* 2018; 50: 915–927
- Meeuwisse WH, Tyreman H, Hagel B et al. A dynamic model of etiology in sport injury: The recursive nature of risk and causation. *Clin J Sport Med* 2007; 17: 215–219
- Bittencourt NFN, Meeuwisse WH, Mendonça LD et al. Complex systems approach for sports injuries: Moving from risk factor identification to injury pattern recognition – narrative review and new concept. *Br J Sports Med* 2016; 50: 1309–1314
- Robertson S. Improving load/injury predictive modelling in sport: The role of data analytics. *J Sci Med Sport* 2014; 18: 25–26
- Combrisson E, Jerbi K. Exceeding chance level by chance: The caveat of theoretical chance levels in brain signal classification and statistical assessment of decoding accuracy. *J Neurosci Methods* 2015; 250: 126–136
- Biran O, Cotton C. Explanation and justification in machine learning: A survey. In: *IJCAI-17 workshop on explainable AI (XAI)*. International Joint Conference on Artificial Intelligence; 2017
- Bellazzi R, Zupan B. Predictive data mining in clinical medicine: Current issues and guidelines. *Int J Med Inform* 2008; 77: 81–97
- Ruddy JD, Shield AJ, Maniar N et al. Predictive modeling of hamstring strain injuries in elite Australian footballers. *Med Sci Sports Exerc* 2018; 50: 906–914
- Dower C, Rafeli A, Weber J, Mohamad R. An enhanced metric of injury risk utilizing Artificial Intelligence. In: *Proceedings of the 13<sup>th</sup> annual MIT SLOAN Sports Analytics Conference*. 2019
- Pasanen K, Rossi MT, Parkkari J et al. Predictors of lower extremity injuries in team sports (PROFITS-study): A study protocol. *BMJ Open Sport Exerc Med* 2015; 1: e000076
- Harriss DJ, MacSween A, Atkinson G. Ethical standards in sport and exercise science research: 2020 update. *Int J Sports Med* 2019; 40: 813–817
- Leppänen M, Pasanen K, Kujala UM et al. Stiff landings are associated with increased ACL injury risk in young female basketball and floorball players. *Am J Sports Med* 2017; 45: 386–393
- Fuller CW, Molloy MG, Bagate C et al. Consensus statement on injury definitions and data collection procedures for studies of injuries in rugby union. *Br J Sports Med* 2007; 41: 328–331
- Boulesteix A-L, Janitza S, Kruppa J et al. Overview of random forest methodology and practical guidance with emphasis on computational biology and bioinformatics. *Wiley Interdiscip Rev. WIREs Data Mining Knowl Discov* 2012; 2: 493–507
- Carey DL, Ong K, Whiteley R et al. Predictive modelling of training loads and injury in Australian football. *Int J Comput Sci Sport* 2018; 17: 49–66
- Breiman L. Random forests. *Mach Learn* 2001; 45: 5–32
- Fawcett T. An introduction to ROC analysis. *Pattern Recognit Lett* 2006; 27: 861–874
- Krstajic D, Buturovic LJ, Leahy DE et al. Cross-validation pitfalls when selecting and assessing regression and classification models. *J Cheminform* 2014; 6: 10–25
- Vanderlei FM, Bastos FN, Tsutsumi GYC et al. Characteristics and contributing factors related to sports injuries in young volleyball players. *BMC Res Notes* 2013; 6: 415
- Jones BH, Bovee MW, Harris JM III et al. Intrinsic risk factors for exercise-related injuries among male and female army trainees. *Am J Sports Med* 1993; 21: 705–710
- Boden BP, Dean GS, Feagin JA et al. Mechanisms of anterior cruciate ligament injury. *Orthopedics* 2000; 23: 573–578
- Yamamoto T. Relationship between hamstring strains and leg muscle strength. A follow-up study of collegiate track and field athletes. *J Sports Med Phys Fitness* 1993; 33: 194–199
- Beynnon BD, Renström PA, Alosa DM et al. Ankle ligament injury risk factors: A prospective study of college athletes. *J Orthop Res* 2001; 19: 213–220
- Woodford-Rogers B, Cyphert L, Denegar CR. Risk factors for anterior cruciate ligament injury in high school and college athletes. *J Athl Train* 1994; 29: 343–346
- Boden BP, Torg JS, Knowles SB et al. Video analysis of anterior cruciate ligament injury: Abnormalities in hip and ankle kinematics. *Am J Sports Med* 2009; 37: 252–259
- Hewett TE, Myer GD, Ford KR et al. Biomechanical measures of neuromuscular control and valgus loading of the knee predict anterior cruciate ligament injury risk in female athletes: A prospective study. *Am J Sports Med* 2005; 33: 492–501