



This is a peer-reviewed, post-print (final draft post-refereeing) version of the following published document and is licensed under Creative Commons: Attribution-Noncommercial-No Derivative Works 4.0 license:

**Oliver, Jon L, Ayala, Francisco, De Ste Croix, Mark B ORCID logo**  
**ORCID: <https://orcid.org/0000-0001-9911-4355>, Lloyd, Rhodri S, Myer, Gregory D and Read, Paul J (2020) Using machine learning to improve our understanding of injury risk and prediction in elite male youth football players. Journal of Science and Medicine in Sport, 23 (11). pp. 1044-1048. doi:10.1016/j.jsams.2020.04.021**

Official URL: [https://www.jsams.org/article/S1440-2440\(19\)31676-7/fulltext](https://www.jsams.org/article/S1440-2440(19)31676-7/fulltext)

DOI: <http://dx.doi.org/10.1016/j.jsams.2020.04.021>

EPrint URI: <https://eprints.glos.ac.uk/id/eprint/8341>

#### **Disclaimer**

The University of Gloucestershire has obtained warranties from all depositors as to their title in the material deposited and as to their right to deposit such material.

The University of Gloucestershire makes no representation or warranties of commercial utility, title, or fitness for a particular purpose or any other warranty, express or implied in respect of any material deposited.

The University of Gloucestershire makes no representation that the use of the materials will not infringe any patent, copyright, trademark or other property or proprietary rights.

The University of Gloucestershire accepts no liability for any infringement of intellectual property rights in any material deposited but will remove such material from public view pending investigation in the event of an allegation of any such infringement.

PLEASE SCROLL DOWN FOR TEXT.

# Using machine learning to improve our understanding of injury risk and prediction in elite male youth football players

**Running head:** Injury risk in youth football

**Authors** Jon L. Oliver<sup>1,2</sup>, Francisco Ayala<sup>3</sup>, Mark B.A. De Ste Croix<sup>4</sup>, Rhodri S. Lloyd<sup>1,2,5</sup>, Greg D. Myer<sup>6</sup> and Paul J. Read<sup>7</sup>

## Affiliations

<sup>1</sup>Youth Physical Development Centre, Cardiff School of Sport and Health Sciences, Cardiff Metropolitan University, Cyncoed Campus, Cardiff Wales, UK.

<sup>2</sup>Sport Performance Research Institute New Zealand (SPRINZ), Auckland University of Technology, Auckland, New Zealand.

<sup>3</sup>Sports Research Centre, Miguel Hernandez University of Elche, Alicante, Spain

<sup>4</sup>School of Physical Education, Faculty of Sport, Health and Social Care, University of Gloucester, United Kingdom

<sup>5</sup>Centre for Sport Science and Human Performance, Waikato Institute of Technology, Hamilton, New Zealand.

<sup>6</sup>Division of Sports Medicine, Cincinnati Children's Hospital, Cincinnati, Ohio, USA

<sup>7</sup>Athlete Health and Performance Research Centre, Aspetar Orthopaedic and Sports Medicine Hospital, Doha, Qatar

## Corresponding Author

Jon L Oliver,  
Youth Physical Development Centre,  
Cardiff Metropolitan University,  
Cyncoed Road,  
Cardiff,  
United Kingdom  
CF23 6XD  
e-mail [joliver@cardiffmet.ac.uk](mailto:joliver@cardiffmet.ac.uk)

Word Count: 3158

## Using machine learning to improve our understanding of injury risk and prediction in elite male youth football players

### Abstract

**Objectives:** The purpose of this study was to examine whether the use of machine learning improved the ability of a neuromuscular screen to identify injury risk factors in elite male youth football players.

**Methods:** 355 elite youth football players aged 10 to 18 years old completed a prospective pre-season neuromuscular screen that included anthropometric measures of size, as well as single leg countermovement jump (SLCMJ), single leg hop for distance (SLHD), 75% hop distance and stick (75% Hop), Y-balance anterior reach and tuck jump assessment. Injury incidence was monitored over one competitive season. Risk profiling was assessed using traditional regression analyses and compared to supervised machine learning algorithms constructed using decision trees.

**Results:** Using continuous data, multivariate logistic analysis identified SLCMJ asymmetry as the 14 sole significant predictor of injury (OR 0.94, 0.92-0.97,  $p < 0.001$  with a specificity of 97.7% and sensitivity of 15.2% giving an AUC of 0.661. The best performing decision tree model provided a specificity of 74.2% and sensitivity of 55.6% with an AUC of 0.663. All variables contributed to the final machine model, with asymmetry in the SLCMJ, 75% Hop and Y-balance, plus tuck jump knee valgus and anthropometrics being the most frequent contributors.

**Conclusions:** Although both statistical methods reported similar accuracy, logistic regression provided very low sensitivity and only identified a single neuromuscular injury risk factor. The machine learning model provided much improved sensitivity to predict injury and identified interactions of asymmetry, knee valgus angle and body size as contributing factors to an injurious profile in youth football players.

**Keywords** Neuromuscular, screen, prospective, binary logistic regression

## 1. INTRODUCTION

Injury rates in elite male youth football can be considered to be high, with recent evidence suggesting that injury rates in this population have increased substantially over the past 15 years due to increased exposure and earlier specialisation.<sup>1,2</sup> Prospective studies have indicated that neuromuscular screening is associated with increased injury risk in professional female and male football players,<sup>3,4</sup> in addition to being associated with anterior cruciate ligament (ACL) and lower extremity injury in elite youth football players.<sup>5,6</sup> Given that neuromuscular control is modifiable,<sup>7</sup> development of a sensitive neuromuscular screen to profile injury risk would be a useful tool to help inform intervention by practitioners.

The ability of a screen to identify injury risk factors and predict injury is often examined using logistic regression. Logistic regression does not manage imbalanced data sets well and tends to only consider the ability of one or a few variables to predict injury, yet it is acknowledged that injury is multifaceted and there will be interactions between and even within risks.<sup>3,8</sup> For example, body size, maturity and neuromuscular control may all interact to influence injury risk in young populations. Consequently, it has been argued that the complexity of injury means a broader statistical approach than logistic regression is needed to better understand relationships between risk factors and predictors of injury.<sup>9</sup>

Machine learning offers a contemporary statistical approach where algorithms have been specifically designed to deal with imbalanced data sets and enable the modelling of interactions between a large number of variables.<sup>3</sup> Contemporary empirical evidence has shown machine learning to provide promising results in the prediction of injury in adult football and handball players,<sup>3,10,11</sup> and high levels of sensitivity in predicting injury in elite youth soccer players from measures of anthropometry and motor performance.<sup>8</sup> To the authors' knowledge, no previous research has provided a direct comparison of the sensitivity and specificity of different statistical approaches in their ability to detect injury risk in athletic populations, with only one existing study employing machine learning to predict injury risk in elite male youth football players.<sup>8</sup> Therefore, the aim of this study was to examine whether the use of machine learning could improve the ability of a neuromuscular screen to predict injury and identify associations between injury risk factors in elite male youth football players.

## 2. MATERIALS AND METHODS

Six professional English Premier League and Championship football clubs volunteered to take part. The clubs initially provided access to  $n = 400$  players, of which  $n = 355$  were tested and prospectively followed for one season, with  $n = 80$  players in the U11s and U12s,  $n = 114$  in the U13s and U14s,  $n = 117$  in the U15s and U16s and  $n = 44$  in the U18s. Participants were included in the study if they were free from illness and injury at the time of testing and regularly involved in training and competitions. Ethical approval was granted by the institutional ethics committee in accordance with the declaration of Helsinki. Parental consent to participate in the study was obtained together with assent from participants.

A prospective cohort study design was used. Following a familiarization session, players were required to attend their respective club's training ground during the pre-season period (July) to undertake a field-based screening battery. Players were then tracked for a period of 10 months (August to June) during the 2014-2015 season to prospectively record all injuries sustained in training and competition. The ability of test variables to identify risk factors and predict injury was compared

using two different statistical methods; univariate and multivariate binary logistic regression analyses versus a selection of popular machine learning methods.

Descriptive variables of chronological age, body mass, stature, leg length, BMI and estimated somatic maturity offset<sup>12</sup> were included along with a battery of neuromuscular control tests; single leg countermovement jump (SLCMJ), single leg hop for distance (SLHD), 75% hop and stick (75%Hop), y-balance anterior reach distance and knee valgus during the tuck jump assessment. Protocols and reliability for all tests have been described elsewhere.<sup>13-16</sup> The SLCMJ and 75%Hop were performed on a force plate (Pasco, Roseville, California, USA) and peak vertical ground reaction force were measured, with all variables reported relative to body weight.<sup>15,17</sup> All tests were performed unilaterally with an asymmetry index calculated for SLCMJ, SLHD, 75%Hop and y-balance.<sup>6</sup>

The procedures for reporting injury occurrence have previously been described elsewhere.<sup>6</sup> Injuries were recorded by medical staff at each football club in accordance with the Premier League's Elite Player Performance Plan. Only non-contact lower-limb injuries were considered for the present study. Injuries were recorded where they resulted from football-related activities and resulted in a player being unable to participate in training or competition for at least 48 hours post incident, not including the day of injury. An injury was classified as non-contact where no clear contact or collision with another player, object or ball occurred. Due to the confounding effects of previous injuries,<sup>18,19</sup> only the first incident experienced by each player during the season was used in the analyses.<sup>4</sup>

Following a traditional statistical approach, a univariate binary logistic regression for each variable was employed. Neuromuscular and anthropometric risk factors that displayed a  $p$  value  $< 0.1$  were considered for further analysis in a multivariate binary logistic regression. Prior to the multivariate analysis the absence of multicollinearity between variables was confirmed using linear regression, with a voluntary inflation factor of  $< 10$  indicating independence between variables. The odds ratio (OR) for each risk factor in the univariate and multivariate analyses were calculated to show the odds of increased injury per unit increase in the independent variable, along with 95% confidence intervals (CI), together with the area under the curve (AUC) with 95% CIs, sensitivity and specificity. This process was repeated for the univariate and multivariate analysis with all variables first analysed as continuous data and secondly as categorical data based on the discretization process described below.

In the contemporary modelling approach, supervised machine learning was employed. Prior to analysis, continuous data were discretized as this can improve the performance of decision trees.<sup>20</sup> Discretization categories for each variable are shown in Table 1. Participants were split into four age groups as per previous research.<sup>6</sup> A single cut-off was applied to asymmetry measures to indicate participants more at risk of injury. Cut-off values for increased injury risk were based on values proposed in previous research; y-balance  $\geq 4$  cm,<sup>16,21</sup> SLCMJ, SLHD and 75%Hop asymmetry all  $\geq 10\%$ .<sup>22</sup> The maturity offset was used to categorise participants as pre ( $> -1$  y), circa ( $-1$  to  $+1$  y) and post ( $> +1$  y) peak height velocity. Knee valgus was categorised as being either absent ( $\leq 0^\circ$ ), minor ( $1-9^\circ$ ), moderate ( $10-20^\circ$ ) or severe ( $> 20^\circ$ ).<sup>23</sup> Following the procedures of Lopez-Valenciano *et al.*,<sup>3</sup> the remaining variables were discretized into four intervals using the unsupervised discretization algorithm available in Weka (Waikato Environment Knowledge Analysis).

**Table 1** Discretization intervals for all variables as applied in the univariate and multivariate binary logistic analysis of categorical data (values for the first group in each variable were used as the reference data for calculation of OR) and applied in all machine learning analyses

Variable	Labels
Age Group	U11-U12, U13-U14, U15-U16 or U18
Body mass (kg)	<43.4, 43.4-53.285, >53.285-65.05 or >65.05
Stature (cm)	<154.025, 154.025-165.55, >165.55-175.85 or >175.85
BMI (kg/m <sup>2</sup> )	<18.145, 18.145-19.565, >19.565-21.545 or >21.545
Maturation offset (y)	>-1, -1 to +1, >+1
Leg length (cm)	<80.55, 80.55-87.25, >87.25-92.08 or >92.08
75% Hop (BW) Left	<2.805, 2.805-3.265, >3.265-3.775 or >3.775
75% Hop (BW) Right	<2.805, 2.805-3.315, >3.315-3.87 or >3.87
75% Hop PVGRF Asym (%)	≤90, >90
SLCMJ (BW) Left	<2.63, 2.63-3.065, >3.065-3.5 or >3.5
SLCMJ (BW) Right	<2.635, 2.635-2.96, >2.96-3.36 or >3.36
SLCMJ PVGRF Asym (%)	≤90, >90
SLHD (leg lengths) Left	<1.525, 1.525-1.675, >1.675-1.825 or >1.825
SLHD (leg lengths) Right	<1.535, 1.535-1.685, >1.685-1.865 or >1.865
SLHD Asym (%)	≤90, >90
TJ Knee Valgus (°) Left	≤0 (none), 1-9 (minor), 10-20 (moderate), >20 (severe)
TJ Knee Valgus (°) Right	≤0 (none), 1-9 (minor), 10-20 (moderate), >20 (severe)
Y-Balance (% leg length) Left	<61.705, 61.705-68.305, >68.305-79.865 or >79.865
Y-Balance (% leg length) Right	<62.375, 62.375-70.675, >70.675-81.31 or >81.31
Y-Balance Asym (cm)	<4, ≥4

Asym = asymmetries; BMI = body mass index; BW = Bodyweights; PVGRF = Peak vertical ground reaction force; SLCMJ = Single leg countermovement jump; SLHD = Single leg hop for distance; TJ = tuck jump;

Three widely used classic decision tree algorithms were chosen as base classifiers; J48 consolidated 123 (J48con), an alternating decision tree (ADT) and a reduces error pruning tree (REPTree), with further 124 algorithms then applied to reduce class imbalance. To address the issue of imbalance and skewed 125 distributions, four resampling, three classic ensemble, three bagging ensemble, three boosting 126 ensemble and five cost-sensitive algorithms were applied to the data. A brief description of each of 127 the techniques employed is provided in supplementary Table S1 and further descriptions of the 128 techniques used are provided by Lopez-Valenciano et al.<sup>3</sup> With all algorithms applied to all base 129 classifiers, a total of 57 models were generated. To allow comparison of the constructed models to a 130 baseline model a ZeroR classifier was also used.

The number of internal classifiers was set at 10 for all ensemble techniques. Thus, each model built by each ensemble technique contained 10 classifying decision trees, each contributed a vote of “yes” or “no” as to whether a participant will get injured. With non-boosting techniques the number of yes/no votes was used to obtain the final prediction, with ≥5 “yes” votes classifying a participant as injured. With boosting techniques each vote was weighted and all votes summed to provide the final prediction, with a summed value >0 classifying a participant as injured. In order to evaluate the performance of the models the data was split into five sets and the five-fold stratified cross validation technique used. For each set, the algorithm was trained with the examples contained in the remaining four sets and then tested with the current set. The AUC was used to evaluate overall accuracy together with the specificity and sensitivity of each model. Cross-validation was used to choose the best performing machine model based on achieving high sensitivity and accuracy.

### 3. RESULTS

Participants had a mean age of  $14.3 \pm 2.1$  y, mass of  $54.3 \pm 13.5$  kg, stature of  $162.4 \pm 14.3$  cm, leg length of  $86.6 \pm 8.2$  cm and maturity offset of  $0.11 \pm 1.93$  y. A total of  $n = 99$  players suffered a first non-contact lower extremity injury during the competitive season. Over three quarters of injuries were classified as either moderate or severe (82%), with the remainder minor or slight. There were a high proportion of strain type injuries (35%), with ligament (17%) and growth/overuse (14%) the most prevalent thereafter.

Results of the univariate analysis for continuous data are presented in Table 2, showing that age, SLCMJ peak force on the right leg, SLCMJ asymmetry and 75% Hop asymmetry were all significantly associated with injury. SLCMJ asymmetry provided an AUC of 0.645 (0.577-0.712) with a sensitivity of 11.1% and specificity of 97.7%. For all other variables, the AUC was  $\leq 0.58$  (95% CIs 0.433-0.645) with a sensitivity of 0% and specificity of 100%. No multicollinearity was present and the multivariate analysis included four variables, providing an AUC of 0.661 (0.596-0.725) with a specificity of 97.7% and sensitivity of 15.2%. In that model only SLCMJ asymmetry provided a significant contribution ( $p < 0.001$ , OR = 0.94, 0.92-0.94), with non-significant contributions from 75% Hop asymmetry ( $p = 0.10$ , OR = 0.98, 0.95-1.00), SLCMJ relative peak force on the right leg ( $p = 0.13$ , OR = 0.72, 0.48-1.10) and age ( $p = 0.36$ , OR = 1.06, 0.94-1.20).

**Table 2** Descriptive statistics and univariate odds ratios from continuous data for all injured and non-injured players

Neuromuscular Risk Factors	Injured Players	Non-injured Players	Odds Ratio (95% CI)	p Value	AUC
Age (y)	14.7 $\pm$ 2.1	14.2 $\pm$ 2.0	1.12 (1.00 - 1.26)	0.05*	0.560
Height (cm)	165.5 $\pm$ 14.1	163.3 $\pm$ 13.6	1.01 (0.99 - 1.03)	0.18	0.542
Mass (kg)	55.7 $\pm$ 14.2	53.7 $\pm$ 13.3	1.01 (0.99 - 1.02)	0.23	0.550
BMI (kg/m <sup>2</sup> )	19.9 $\pm$ 2.4	19.8 $\pm$ 2.4	1.02 (0.92 - 1.12)	0.63	0.515
Leg Length (cm)	86.9 $\pm$ 7.6	85.9 $\pm$ 8.4	1.02 (0.99 - 1.05)	0.25	0.534
Maturity-Offset	0.3 $\pm$ 1.9	0.1 $\pm$ 1.9	0.95 (0.84 - 1.07)	0.37	0.531
75% Hop L PVGRF (BW)	3.37 $\pm$ 0.65	3.25 $\pm$ 0.69	1.31 (0.93 - 1.83)	0.12	0.554
75% Hop R PVGRF (BW)	3.43 $\pm$ 0.82	3.34 $\pm$ 0.74	1.16 (0.86 - 1.58)	0.32	0.520
75% Hop Asym (%)	86.2 $\pm$ 9.3	88.2 $\pm$ 8.1	0.97 (0.95 - 1.00)	0.05*	0.557
SLCMJ L PVGRF (BW)	3.07 $\pm$ 0.65	3.11 $\pm$ 0.64	0.91 (0.63 - 1.32)	0.64	0.502
SLCMJ R PVGRF (BW)	2.96 $\pm$ 0.61	3.10 $\pm$ 0.60	0.66 (0.44 - 0.99)	0.05*	0.577
SLCMJ PVGRF Asym (%)	82.9 $\pm$ 9.7	87.6 $\pm$ 7.8	0.94 (0.91 - 0.97)	<0.001**	0.645
SLHD L (% leg length)	1.72 $\pm$ 0.3	1.69 $\pm$ 0.3	1.37 (0.63 - 2.99)	0.42	0.521
SLHD R (% leg length)	1.74 $\pm$ 0.3	1.71 $\pm$ 0.3	1.34 (0.62 - 2.90)	0.45	0.523
SLHD Asym (%)	93.1 $\pm$ 5.7	94.0 $\pm$ 5.0	0.97 (0.93 - 1.01)	0.19	0.544
TJ Knee Valgus L	1.07 $\pm$ 0.9	1.19 $\pm$ 0.9	0.85 (0.66 - 1.12)	0.26	0.545
TJ Knee Valgus R	1.43 $\pm$ 0.9	1.34 $\pm$ 0.9	1.12 (0.86 - 1.46)	0.39	0.523
Y-B (% leg length) L	70.6 $\pm$ 13.3	71.1 $\pm$ 14.5	0.99 (0.98 - 1.01)	0.79	0.505
Y-B (% leg length) R	73.2 $\pm$ 13.7	71.9 $\pm$ 15.0	1.00 (0.99 - 1.02)	0.47	0.534
Y-B Asym (%)	94.0 $\pm$ 4.8	94.0 $\pm$ 5.0	1.00 (0.96 - 1.05)	0.92	0.500

\* Significant at the level of  $p < .05$

\*\*Significant at the level of  $p < .001$

BMI = Body mass index; Asym = asymmetry; BW = body weight; SLCMJ = single leg countermovement jump; SLHD = single leg hop for distance; TJ = Tuck Jump; PVGRF = peak vertical ground reaction force; Y-B = y-balance; 75% Hop = 75% horizontal hop and stick; R = right; L = left

Predictive ability was similar when logistic regression was performed on categorical data; in the univariate analysis all variables reported an AUC  $\leq 0.57$  (95% CIs 0.435-0.641) with a sensitivity of 0% and specificity of 100%. With the absence of collinearity, a multivariate analysis including height, 75% Hop asymmetry, SLCMJ on the right leg and asymmetry, Y-balance on both the left and right leg and tuck jump knee valgus on the left leg provided a prediction with a specificity of 94.5% and sensitivity of 11.1% and an AUC of 0.687 (0.627-0.747). The multivariate analysis is available in Supplementary Table S2 and shows SLCMJ PVGRF on the right leg as the only significant predictor across groups (all OR 0.49, 0.25 – 0.98).

With machine learning, the baseline ZeroR classifier achieved an AUC of 0.494, specificity of 100% and sensitivity of 0%. Supplementary Tables S3, S4 and S5 show the performance of the different decision trees for the resampling, ensemble and cost-sensitive machine learning techniques respectively, nearly all of which have greater accuracy and sensitivity than the baseline model. The bagging ensemble method with a J48con decision tree as a base classifier and a 1:1 cost sensitive

learning matrix was chosen as the best performing decision tree model. Cross-validation showed an AUC of 0.663 (0.550-0.776) with the model correctly classifying 74.2% of non-injured and 55.6% of injured players. All classifiers from the final model are available in supplementary Figures S1-S10, with Figure S1 showing a decision pathway for an example player. Table 3 shows the frequency with which each of the 20 risk factors appeared across the 10 classifiers in the final model. SLCMJ asymmetry appeared in all classifiers, with SLHD asymmetry, hop and stick (75% Hop) asymmetry and knee valgus on the left leg identified as other frequently included neuromuscular risks ( $\geq 7/10$ ). A number of descriptive measures also appeared frequently in the final model, including age, body mass, stature and leg length.

**Table 3** Number of classifiers (out of 10) in which each risk factor featured in the final machine learning model

<b>Risk Factor</b>	<b>N° of Classifiers</b>
SLCMJ PVGRF Asym	10
Body mass	8
Leg length	8
Stature	8
Age group	7
75% Hop PVGRF Asym	7
SLHD Asym	7
TJ Knee Valgus Left	7
SLHD Left	6
SLCMJ PVGRF Left	5
SLCMJ PVGRF Right	5
75% Hop Left	4
Maturation offset	4
SLHD Right	4
Y-Balance Asym	4
Y-Balance Left	4
BMI	3
75% Hop Right	3
TJ Knee Valgus Right	3
Y-Balance Right	2

Asym = asymmetry; BMI = body mass index; PVGRF = peak vertical ground reaction force SLCMJ = single leg countermovement jump; SLHD = single leg hop for distance; TJ = tuck jump

#### 4. DISCUSSION

The aim of the present study was to examine whether the use of machine learning improved the ability to identify injury risk factors and predict injury in a cohort of elite male youth football players. Machine learning did not improve the overall accuracy of injury prediction. However, logistic regression was heavily biased towards the majority class of non-injured players and provided poor sensitivity, whereas machine learning provided a more balanced predictive model with sensitivity to predict injury improving more than 3.5 -fold. Whether using continuous or categorical data, multivariate logistic regression only identified a single significant predictor of injury and improved sensitivity in the machine learning model may reflect a better ability of that model to consider interactions between risk factors. All variables appeared multiple times in the final machine learning model, suggesting the importance of interactions between asymmetry, movement control and body size as injury risk factors in elite youth football.

It has been suggested that using logistic regression does not control well for imbalanced data sets when predicting injury.<sup>3</sup> This seems to be the case in the present study with all univariate analyses showing perfect to near perfect specificity but low to zero sensitivity, indicating a good ability to identify individuals who did not get injured but not those who did get injured. Similarly, multivariate logistic regression using both continuous data and categorical data achieved low levels of sensitivity



( $\leq 15.2\%$ ) with both approaches only identifying one significant predictor of injury. The absence of other significant contributing variables suggests that multivariate logistic regression may not be proficient at quantifying interactions between risk factors.

The identification of relatively few significant predictor variables with logistic regression may be partly due to pooling of data for elite youth players across a broad range of age groups comprising different stages of maturation. In U14 to U16 age groups younger players experience more overuse injuries,<sup>24</sup> while in U11 to U14 late maturing players experience more overuse injuries.<sup>25</sup> Using a similar screen to the present study, Read et al.<sup>6</sup> recently reported that maturity offset was the only variable significantly associated with injury in U13-U14 y old players (OR = 0.58), while heightened SLCMJ asymmetry (OR = 0.90) and lower relative SLCMJ peak force (OR = 0.36) were significantly associated with injury in the U11-12 and U15-16 age categories respectively. Given the complex interaction of growth, maturity timing and tempo, and injury the pooling of data across players of varying age and maturity in the present study may have reduced the ability of logistic regression to successfully identify participants who experienced an injury.

The AUC indicated that machine learning models also had poor overall accuracy to detect injury. This is similar to a previous study examining hamstring strain in elite Australian footballers, which reported that both logistic regression and machine learning achieved low overall accuracy (AUC < 0.60).<sup>26</sup> A recent review of clinical diagnostic tools also reported no advantage of machine learning over logistic regression on overall accuracy.<sup>27</sup> However, researchers should not focus solely on the AUC, but also consider the need for higher sensitivity,<sup>28</sup> as identifying players with an increased risk of injury should be a priority. With the imbalanced data set used in the current study, it appears relatively easy to construct models with high specificity but more difficult to create models with high sensitivity. With the chosen machine learning model, the sensitivity was improved over 3.5 fold, with the model identifying all variables as contributing to injury risk.

Similar to the recent work of Rommers et al.<sup>8</sup> but with lower sensitivity, our machine learning model identified measures of size as important predictors of injury in youth soccer, although our model also noted asymmetry and valgus often contribute to an injurious profile. This may be important as young football players are known to exhibit lower limb asymmetries and reduced frontal plane knee control<sup>23,29</sup> and these are likely to be modifiable risk factors. Our modelling process used information entropy to identify attributes that provided the greatest normalized gain in predicting injury, with subsequent pruning of the decision tree to reduce complexity and over-fitting, remove noisy data and improve predictions. This means that classifiers from the model are relatively straight forward to follow. For our sample population results would suggest that interventions that consider size, maturity, asymmetry and movement control may be useful to reduce injury. Where an individual player is identified at risk of injury practitioners should further explore results to identify which decision tree classifiers and features contributed to that outcome, which neuromuscular risk factors within those classifiers could be modified to reduce risk and which non-modifiable factors may need to be managed (e.g. by reducing exposure).

In the logistic regression analysis data imbalance was not adjusted and results were validated in the population with which the predictive equations were generated, which will create a bias and inflate the overall accuracy. The machine learning adopted a more robust stratified cross-validation approach, making accuracy more difficult to achieve. The statistical approaches used were purposely chosen to reflect the manner in which they are typically applied to injury profiling and prediction. To maintain sample size and reduce overfitting, the present study examined all non-contact lower limb

injuries and did not focus on a particular injury type or severity classification. A focus on a single injury type (e.g. ACL) or classification (e.g. muscle injury) has been shown to provide better accuracy in previous research using logistic regression<sup>5</sup> and machine learning,<sup>3,10</sup> but this approach would have led to data reduction and greater data imbalance.

## **5. CONCLUSION**

The ability to predict injury and our understanding of the factors that contribute to injury risk are influenced by the statistical approaches used to analyse prospective data, which may then influence practice. Most likely due to the complex nature of injury occurrence and prediction, achieving a high level of overall accuracy may be difficult with both logistic regression (whether using continuous or categorical data) and machine learning approaches. If more importance is placed on sensitivity (rather than specificity) then machine learning may offer a promising method to predict injury, while also providing a deeper understanding of the interaction between variables that contribute to injury risk. In the cohort examined, machine learning suggested that asymmetry, knee valgus angle, age and size all contribute to injury risk in elite male youth football players. Given that movement mechanics and asymmetry are modifiable qualities, these findings may help guide injury prevention practice and future research.

## **6. PRACTICAL APPLICATIONS**

- Machine learning improved the sensitivity of injury prediction more than 3.5-fold compared to multivariate logistic regression analyses.
- Asymmetry, movement control, maturity and size all interact to influence injury risk in elite male youth football players, supporting the need to screen for a variety of risk factors and consider results collectively.
- Given that movement control (e.g. knee valgus) and asymmetry are modifiable neuromuscular risk factors, at risk players may benefit from interventions that target deficits in these qualities.

## **ACKNOWLEDGEMENTS**

The authors would like to thank the clubs and individuals who participated in the research. There was no financial assistance for this research.

## REFERENCES

1. Read PJ, Oliver JL, De Ste Croix MB, Myer GD, Lloyd RS. The scientific foundations and associated injury risks of early soccer specialisation. *J Sports Sci.* 2016;34(24):2295-2302.
2. Read PJ, Oliver JL, De Ste Croix MBA, Myer GD, Lloyd RS. An audit of injuries in six english professional soccer academies. *J Sports Sci.* 2018;36(13):1542-1548.
3. Lopez-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(5):915-927.
4. Nilstad A, Andersen TE, Bahr R, Holme I, Steffen K. Risk factors for lower extremity injuries in elite female soccer players. *Am J Sports Med.* 2014;42(4):940-948.
5. Padua DA, DiStefano LJ, Beutler AI, de la Motte SJ, DiStefano MJ, Marshall SW. The Landing Error Scoring System as a Screening Tool for an Anterior Cruciate Ligament Injury-Prevention Program in Elite-Youth Soccer Athletes. *J Athl Train.* 2015;50(6):589-595.
6. Read PJ, Oliver JL, De Ste Croix MBA, Myer GD, Lloyd RS. A prospective investigation to evaluate risk factors for lower extremity injury risk in male youth soccer players. *Scand J Med Sci Sports.* 2018;28(3):1244-1251.
7. Sugimoto D, Myer GD, Barber Foss KD, Pepin MJ, Micheli LJ, Hewett TE. Critical components of neuromuscular training to reduce ACL injury risk in female athletes: meta-regression analysis. *Br J Sports Med.* 2016;50(20):1259-1266.
8. Rommers N, Rossler R, Verhagen E, et al. A Machine Learning Approach to Assess Injury Risk in Elite Youth Football Players. *Med Sci Sports Exerc.* 2020.
9. Bittencourt NFN, Meeuwisse WH, Mendonca LD, Nettel-Aguirre A, Ocarino JM, Fonseca ST. 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(21):1309-1314.
10. Ayala F, Lopez-Valenciano A, Gamez Martin JA, et al. A Preventive Model for Hamstring Injuries in Professional Soccer: Learning Algorithms. *Int J Sports Med.* 2019;40(5):344-353.
11. Rossi A, Pappalardo L, Cintia P, Iaia FM, Fernandez J, Medina D. Effective injury forecasting in soccer with GPS training data and machine learning. *PLoS One.* 2018;13(7):e0201264.
12. Mirwald RL, Baxter-Jones AD, Bailey DA, Beunen GP. An assessment of maturity from anthropometric measurements. *Med Sci Sports Exerc.* 2002;34(4):689-694.
13. Read PJ, Oliver JL, Croix MB, Myer GD, Lloyd RS. Consistency of Field-Based Measures of Neuromuscular Control Using Force-Plate Diagnostics in Elite Male Youth Soccer Players. *J Strength Cond Res.* 2016;30(12):3304-3311.
14. Read PJ, Oliver JL, de Ste Croix MB, Myer GD, Lloyd RS. Reliability of the Tuck Jump Injury Risk Screening Assessment in Elite Male Youth Soccer Players. *J Strength Cond Res.* 2016;30(6):1510-1516.
15. Read PJ, Oliver JL, Myer GD, De Ste Croix MBA, Belshaw A, Lloyd RS. Altered landing mechanics are shown by male youth soccer players at different stages of maturation. *Phys Ther Sport.* 2018;33:48-53.
16. Smith CA, Chimera NJ, Warren M. Association of y balance test reach asymmetry and injury in division I athletes. *Med Sci Sports Exerc.* 2015;47(1):136-141.
17. Read PJ, Oliver JL, De Ste Croix MBA, Myer GD, Lloyd RS. Hopping and Landing Performance in Male Youth Soccer Players: Effects of Age and Maturation. *Int J Sports Med.* 2017;38(12):902-908.

18. Arnason A, Sigurdsson SB, Gudmundsson A, Holme I, Engebretsen L, Bahr R. Risk factors for injuries in football. *Am J Sports Med.* 2004;32(1 Suppl):5S-16S.
19. Kucera KL, Marshall SW, Kirkendall DT, Marchak PM, Garrett WE, Jr. Injury history as a risk factor for incident injury in youth soccer. *Br J Sports Med.* 2005;39(7):462.
20. Hacibeyoglu M, Arslan A, Kahramanli S. Improving classification accuracy with discretization on data sets including continuous valued features. *International Journal of Computer, Electrical, Automation, Control and Information Engineering.* 356 2011;4(6):623-626.
21. Plisky PJ, Rauh MJ, Kaminski TW, Underwood FB. Star Excursion Balance Test as a predictor of lower extremity injury in high school basketball players. *J Orthop Sports Phys Ther.* 2006;36(12):911-919.
22. Kyritsis P, Bahr R, Landreau P, Miladi R, Witvrouw E. Likelihood of ACL graft rupture: not meeting six clinical discharge criteria before return to sport is associated with a four times greater risk of rupture. *Br J Sports Med.* 2016;50(15):946-951.
23. Read PJ, Oliver JL, De Ste Croix MBA, Myer GD, Lloyd RS. Landing Kinematics in Elite Male Youth Soccer Players of Different Chronologic Ages and Stages of Maturation. *J Athl Train.* 2018;53(4):372-378.
24. Le Gall F, Carling C, Reilly T, Vandewalle H, Church J, Rochcongar P. Incidence of injuries in elite French youth soccer players: a 10-season study. *Am J Sports Med.* 2006;34(6):928-938.
25. van der Sluis A, Elferink-Gemser MT, Brink MS, Visscher C. Importance of Peak Height Velocity Timing in Terms of Injuries in Talented Soccer Players. *Int J Sports Med.* 2015;36(4):327-332.
26. Ruddy JD, Shield AJ, Maniar N, et al. Predictive Modeling of Hamstring Strain Injuries in Elite Australian Footballers. *Med Sci Sports Exerc.* 2018;50(5):906-914.
27. Christodoulou E, Ma J, Collins GS, Steyerberg EW, Verbakel JY, van Calster B. A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models. *J Clin Epidemiol.* 2019.
28. Florkowski CM. Sensitivity, specificity, receiver-operating characteristic (ROC) curves and likelihood ratios: communicating the performance of diagnostic tests. *Clin Biochem Rev.* 2008;29 Suppl 1:S83-87.
29. Atkins SJ, Bentley I, Hurst HT, Sinclair JK, Hesketh C. The Presence of Bilateral Imbalance of the Lower Limbs in Elite Youth Soccer Players of Different Ages. *J Strength Cond Res.* 2016;30(4):1007-1013.

## Supplementary Tables

**Table S1** Brief descriptions of the resampling, ensemble and cost-sensitive algorithms applied to the base classifier decision trees

Algorithm	Description
<b>Resampling</b>	
Smote I	Synthetic minority oversampling technique is an oversampling method, whose main idea is to create new minority class examples by interpolating several minority class instances that lie together. SMOTE creates instances by randomly selecting one (or more de- pending on the oversampling ratio) of the k nearest neighbors (kNN) of a minority class instance and the generation of the new instance values from a random interpolation of both instances. In this study, each decision tree applied on data set previously pre-processed with Smote.  They were considered using the k-3 and k-5 nearest neighbour of minority class instances using the Euclidean distance
Smote II	
Random Oversampling (ROS)	Each decision tree applied on the data set previously pre-processed with random over sampling technique, a filter that duplicates some random minority instances until the total amount of minority instances reaches the percentage given.
Random Undersampling (RUS)	Each decision tree applied on the data set previously pre-processed with random under sampling technique, a filter that eliminates some random majority instances until the total amount of majority instances reaches the percentage given
<b>Ensemble</b>	
Adaboost	Classic AdaBoost, without using confidences.  AdaBoost uses the whole data-set to train each classifier serially, but after each round, it gives more focus to difficult instances, with the goal of correctly classifying examples in the next iteration that were incorrectly classified during the current iteration. Hence, it gives more focus to examples that are harder to classify, the quantity of focus is measured by a weight, which initially is equal for all instances. After each iteration, the weights of misclassified instances are increased while the weights of correctly classified instances are decreased. Furthermore, another weight is assigned to each individual classifier depending on its overall accuracy which is then used in the test phase; more confidence is given to more accurate classifiers. Finally, when a new instance is

	submitted, each classifier gives a weighted vote, and the class label is selected by majority.
Adaboost – M1	Multi-class AdaBoost, slightly different weight update.
Bagging	<p>Classic Bagging, resampling with replacement, bag size equal to original data set size.</p> <p>It consists in training different classifiers with bootstrapped replicas of the original training dataset. That is, a new data-set is formed to train each classifier by randomly drawing (with replacement) instances from the original data-set (usually, maintaining the original data-set size). Hence, diversity is obtained with the resampling procedure by the usage of different data subsets. Finally, when an unknown instance is presented to each individual classifier, a majority or weighted vote is used to infer the class.</p>
Smoteboost	<p>AdaBoost.M2 with Smote in each iteration.</p> <p>The weights of the new instances are proportional to the total number of instances in the new data-set. Hence, their weights are always the same (in all iterations and for all new instances), whereas original data-set's instances weights are normalized in such a way that they form a distribution with the new instances. After training a classifier, the weights of the original data-set instances are updated; then another sampling phase is applied (again, modifying the weight distribution). The repetition of this process also brings along more diversity in the training data, which generally benefits the ensemble learning.</p>
RUSboost	<p>AdaBoost.M2 with random undersampling in each iteration.</p> <p>RUSboost performs similarly to SmoteBoost, but it removes instances from the majority class by random undersampling the dataset in each iteration. It is not necessary to assign new weights to the instances. It is enough with simply normalizing the weights of the remaining instances in the new dataset with respect to their total sum of weights.</p>
ROSboost	<p>AdaBoost.M2 with random oversampling in each iteration.</p> <p>Similar to RUSboost, but it creates instances from the minority class by random oversampling the dataset in each iteration.</p>
Overbagging	<p>Bagging with oversampling of the minority class.</p> <p>Instead of performing a random sampling of the whole dataset, an oversampling process can be carried out before training each classifier.</p>
Underbagging	Bagging with undersampling of the majority class.

	On the contrary to Overbagging, Underbagging procedure uses undersampling instead of oversampling.
Smotebagging	Bagging where each bag's Smote quantity varies
<b>Cost-sensitive (CS)</b>	
Meta-Cost	Both consider the variable cost of a misclassification with respect to the different classes. The final cost matrix set-up was based on the best performance reported after testing all the possibilities.
CS-classifier	
CS-Adaboost-M1	Adaboost – M1 with an asymmetric classification cost matrix in the base classifier
CS-Adaboost	Adaboost with an asymmetric classification cost matrix in the base classifier
CS-Bagging	Bagging with an asymmetric classification cost matrix in the base classifier

**Table S2** Multivariate analysis with odds ratios from categorical data. Groups were ordered in size from smallest to largest where there were four groups (see Table 1) with group 1 acting as the reference group. For asymmetry the group with no asymmetry acted as the reference.

<b>Neuromuscular Risk Factor</b>	<b>Beta</b>	<b><i>p</i> Value</b>	<b>OR (95% CI)</b>
<b>Constant</b>	-0.37		
<b>Height</b>			
Group 1		Reference	1
Group 2	0.50	0.18	1.64 (0.79 – 3.41)
Group 3	-0.15	0.71	0.86 (0.40 – 1.85)
Group 4	0.21	0.56	1.24 (0.59 – 2.56)
<b>Y-Balance Right</b>			
Group 1		Reference	1
Group 2	0.34	0.60	1.41 (0.40 – 4.95)
Group 3	0.44	0.38	1.55 (0.58 – 4.13)
Group 4	0.25	0.56	1.29 (0.56 – 2.97)
<b>SLCMJ Right</b>			
Group 1		Reference	1
Group 2	-0.71	0.04*	0.49 (0.25 – 0.97)
Group 3	-0.71	0.04*	0.49 (0.25 – 0.98)
Group 4	-0.71	0.04*	0.49 (0.25 – 0.98)
<b>TJ Valgus Left</b>			
Group 1		Reference	1
Group 2	0.48	0.36	1.62 (0.58 – 4.56)
Group 3	-0.68	0.05	0.51 (0.26 – 1.01)
Group 4	-0.14	0.65	0.87 (0.48 – 1.58)
<b>Y-Balance Left</b>			
Group 1		Reference	1
Group 2	-0.44	0.49	0.64 (0.18 – 2.26)
Group 3	0.31	0.52	1.37 (0.52 – 3.58)
Group 4	-0.18	0.68	0.84 (0.36 – 1.93)
<b>75%Hop Asymmetry</b>			
Group 1		Reference	1
Group 2	-0.15	0.55	0.86 (0.52 – 1.42)
<b>SLCMJ Asymmetry</b>			
Group 1		Reference	1
Group 2	-0.23	0.39	0.80 (0.48 – 1.33)

\*Significantly different to group 1 ( $p < 0.05$ )

75%Hop = 75% Hop and stick; SLCMJ = single leg countermovement jump;

TJ = tuck jump



**Table S3** Average AUC, sensitivity and specificity for all decision trees in isolation and after resampling.

		Technique	AUC	Sensitivity (%)	Specificity (%)
Base Classifiers		ADTree	0.613	30.3	84.4
		J48 <sub>CON</sub>	<b>0.626</b>	<b>66.7</b>	54.7
		REPTree	0.494	7.1	<b>90.2</b>
Smote I (k = 3)					
Oversampling techniques		ADTree	<b>0.604</b>	<b>45.5</b>	<b>75.4</b>
		J48 <sub>CON</sub>	0.571	37.4	68.4
		REPTree	0.547	43.4	69.5
		Smote II (k = 5)			
		ADTree	<b>0.601</b>	40.8	<b>80.1</b>
		J48 <sub>CON</sub>	0.583	<b>42.4</b>	62.9
		REPTree	0.573	40.4	77
		Random Oversampling			
		ADTree	<b>0.623</b>	36.4	77.7
		J48 <sub>CON</sub>	0.584	<b>49.5</b>	69.1
		REPTree	0.608	38.4	<b>78.5</b>
Random Undersampling					
Undersampling techniques		ADTree	<b>0.635</b>	41.4	<b>77</b>
		J48 <sub>CON</sub>	0.616	<b>63.6</b>	52.3
		REPTree	0.61	40.4	<b>77</b>

**Table S4** Average AUC, sensitivity and specificity for the machine learning ensembles techniques

	Technique	AUC	Sensitivity (%)	Specificity (%)
Classic Ensemble	Adaboost			
	ADTree	0.547	23.2	<b>84.8</b>
	J48 <sub>CON</sub>	<b>0.574</b>	<b>31.3</b>	78.5
	REPTree	0.564	29.3	76.6
	Adaboost-M1			
	ADTree	<b>0.592</b>	32.3	<b>82.8</b>
	J48 <sub>CON</sub>	0.552	<b>36.4</b>	71.5
	REPTree	0.589	29.3	78.5
	Bagging			
	ADTree	<b>0.633</b>	27.3	91.4
	J48 <sub>CON</sub>	0.630	<b>46.5</b>	73.4
	REPTree	0.559	9.1	<b>94.1</b>
Boosting Ensembles	SmoteBoost			
	ADTree*	<b>0.621</b>	<b>43.4</b>	72.3
	J48 <sub>CON</sub> *	0.583	34.3	<b>76.2</b>
	REPTree*	0.586	37.4	74.6
	ROSBost			
	ADTree	<b>0.603</b>	27.3	<b>83.6</b>
	J48 <sub>CON</sub>	0.593	<b>32.3</b>	77
	REPTree	0.598	<b>32.3</b>	75
	RUSBoost			
	ADTree	0.637	<b>48.5</b>	76.2
	J48 <sub>CON</sub>	<b>0.644</b>	40.4	77
	REPTree	0.583	10.1	<b>95.7</b>
Bagging Esembles	OverBagging			
	ADTree	0.657	35.4	<b>85.2</b>
	J48 <sub>CON</sub>	<b>0.665</b>	<b>43.4</b>	80.1
	REPTree	0.636	24.2	86.7
	UnderBagging			
	ADTree	<b>0.663</b>	42.4	80.9
	J48 <sub>CON</sub>	0.651	<b>52.5</b>	73
	REPTree	0.629	30.3	<b>82.8</b>
	SmoteBagging			

ADTree*	0.630	<b>42.4</b>	80.5
J48 <sub>CON</sub> <sup>T</sup>	<b>0.657</b>	39.4	84.4
REPTree <sup>T</sup>	0.609	29.3	<b>87.1</b>

\*: Smote with  $k = 3$  nearest neighbours of minority class instances; <sup>T</sup>: Smote with  $k = 5$  nearest neighbours of minority class instances

**Table S5** Average AUC, sensitivity and specificity for the cost-sensitive learning and class-balanced ensembles with a cost-sensitive classifier techniques (grey shaded area shows the selected best-performing model)

	Technique	AUC	Sensitivity (%)	Specificity (%)	Cost matrix
Cost- Sensitive Classifiers	MetaCost				
	ADTree	<b>0.635</b>	<b>59.6</b>	65.2	$\begin{Bmatrix} 0 & 2 \\ 1 & 0 \end{Bmatrix}$
	J48 <sub>CON</sub>	<b>0.635</b>	57.6	<b>71.9</b>	$\begin{Bmatrix} 0 & 1 \\ 1 & 0 \end{Bmatrix}$
	REPTree	0.569	60.6	60.9	$\begin{Bmatrix} 0 & 3 \\ 1 & 0 \end{Bmatrix}$
	CS-Classifier				
	ADTree	0.609	46.5	<b>72.7</b>	$\begin{Bmatrix} 0 & 2 \\ 1 & 0 \end{Bmatrix}$
	J48 <sub>CON</sub>	<b>0.626</b>	<b>66.7</b>	54.7	$\begin{Bmatrix} 0 & 1 \\ 1 & 0 \end{Bmatrix}$
	REPTree	0.521	34.3	71.5	$\begin{Bmatrix} 0 & 2 \\ 1 & 0 \end{Bmatrix}$
Classic ensembles with a cost-sensitive classifier	CS-Adaboost-M1				
	ADTree	<b>0.610</b>	37.4	76.7	$\begin{Bmatrix} 0 & 3 \\ 2 & 0 \end{Bmatrix}$
	J48 <sub>CON</sub>	0.608	<b>43.4</b>	73.4	$\begin{Bmatrix} 0 & 2 \\ 1 & 0 \end{Bmatrix}$
	REPTree	0.591	30.3	<b>80.1</b>	$\begin{Bmatrix} 0 & 3 \\ 2 & 0 \end{Bmatrix}$
	CS-Adaboost				
	ADTree	<b>0.635</b>	30.3	<b>85.9</b>	$\begin{Bmatrix} 0 & 1 \\ 1 & 0 \end{Bmatrix}$
	J48 <sub>CON</sub>	0.623	29.3	81.3	$\begin{Bmatrix} 0 & 1 \\ 1 & 0 \end{Bmatrix}$
	REPTree	0.486	<b>31.3</b>	71.1	$\begin{Bmatrix} 0 & 1 \\ 1 & 0 \end{Bmatrix}$
	CS-Bagging				
	ADTree	0.638	50.5	75.8	$\begin{Bmatrix} 0 & 2 \\ 1 & 0 \end{Bmatrix}$
	<b>J48<sub>CON</sub></b>	<b>0.663</b>	<b>55.6</b>	74.2	$\begin{Bmatrix} 0 & 1 \\ 1 & 0 \end{Bmatrix}$
	REPTree	0.653	36.4	<b>78.9</b>	$\begin{Bmatrix} 0 & 3 \\ 1 & 0 \end{Bmatrix}$



Figure S 2

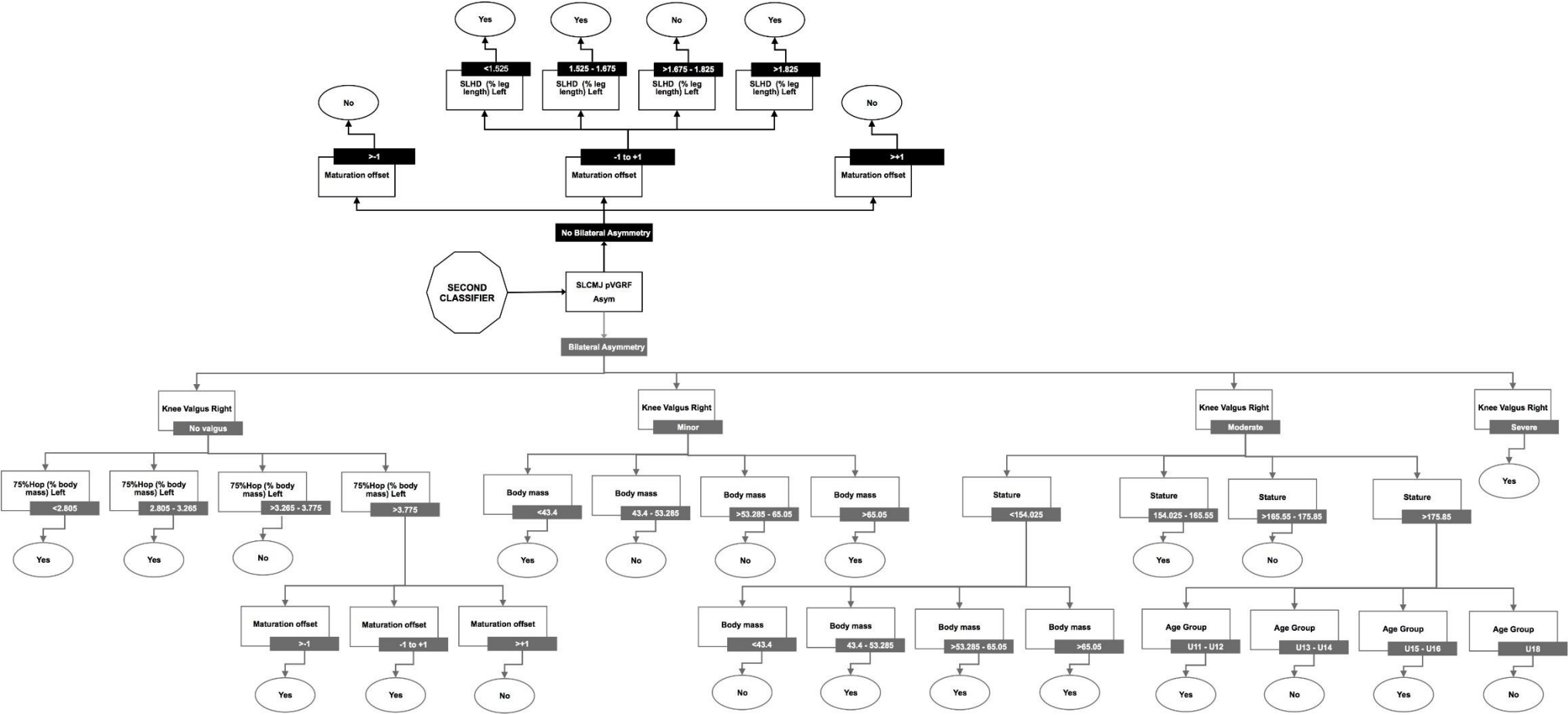
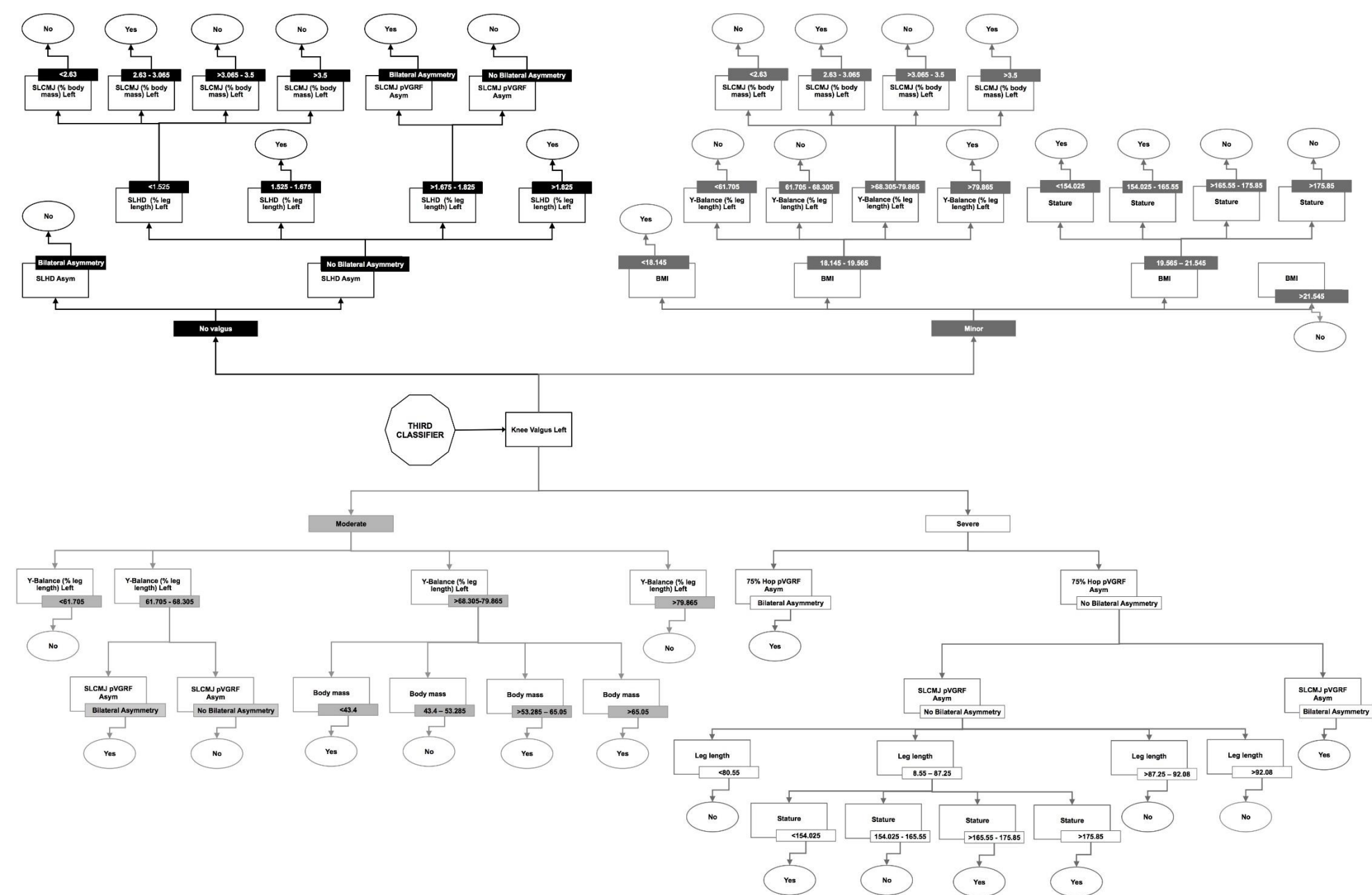


Figure S 3



**Figure S 4**

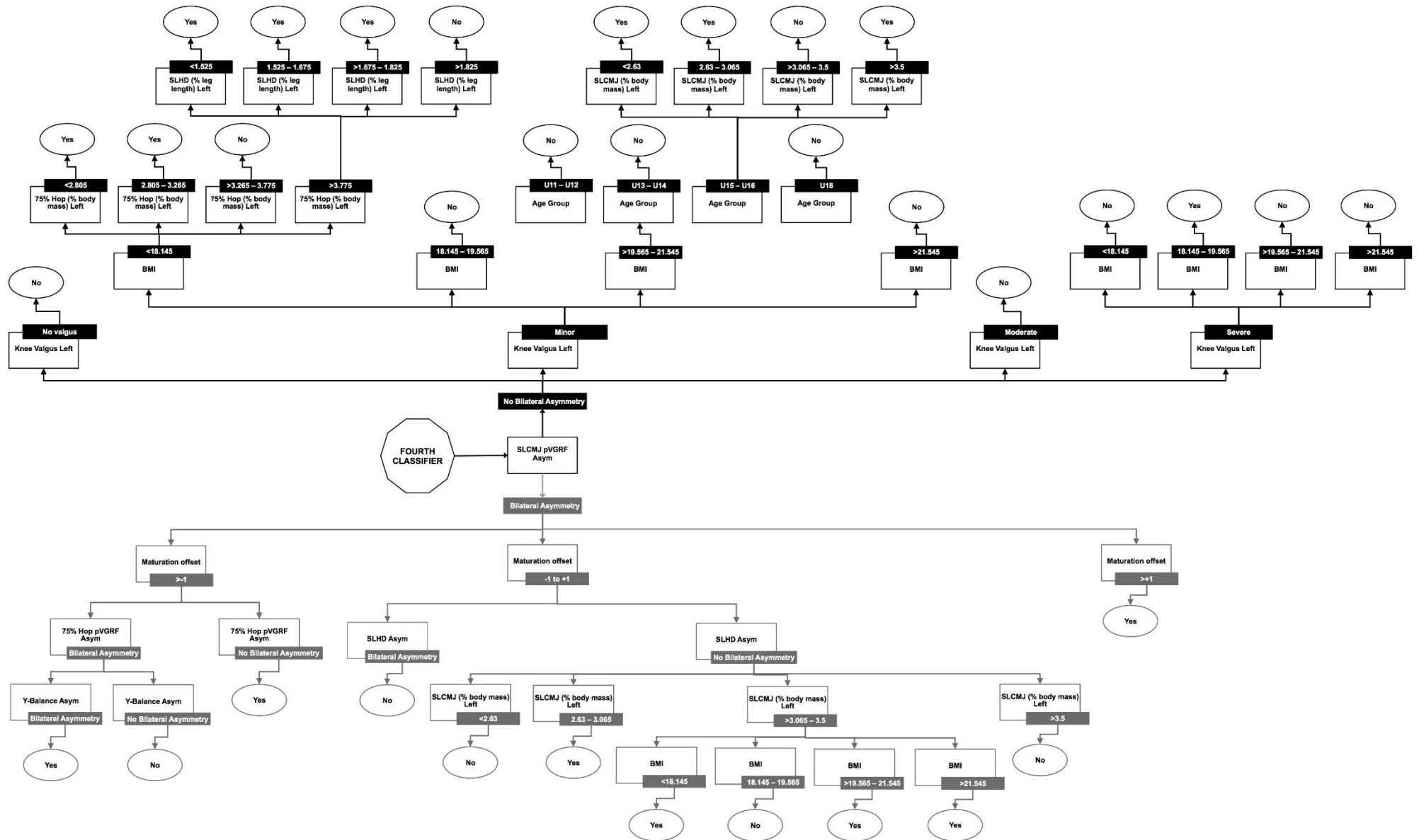




Figure S 5

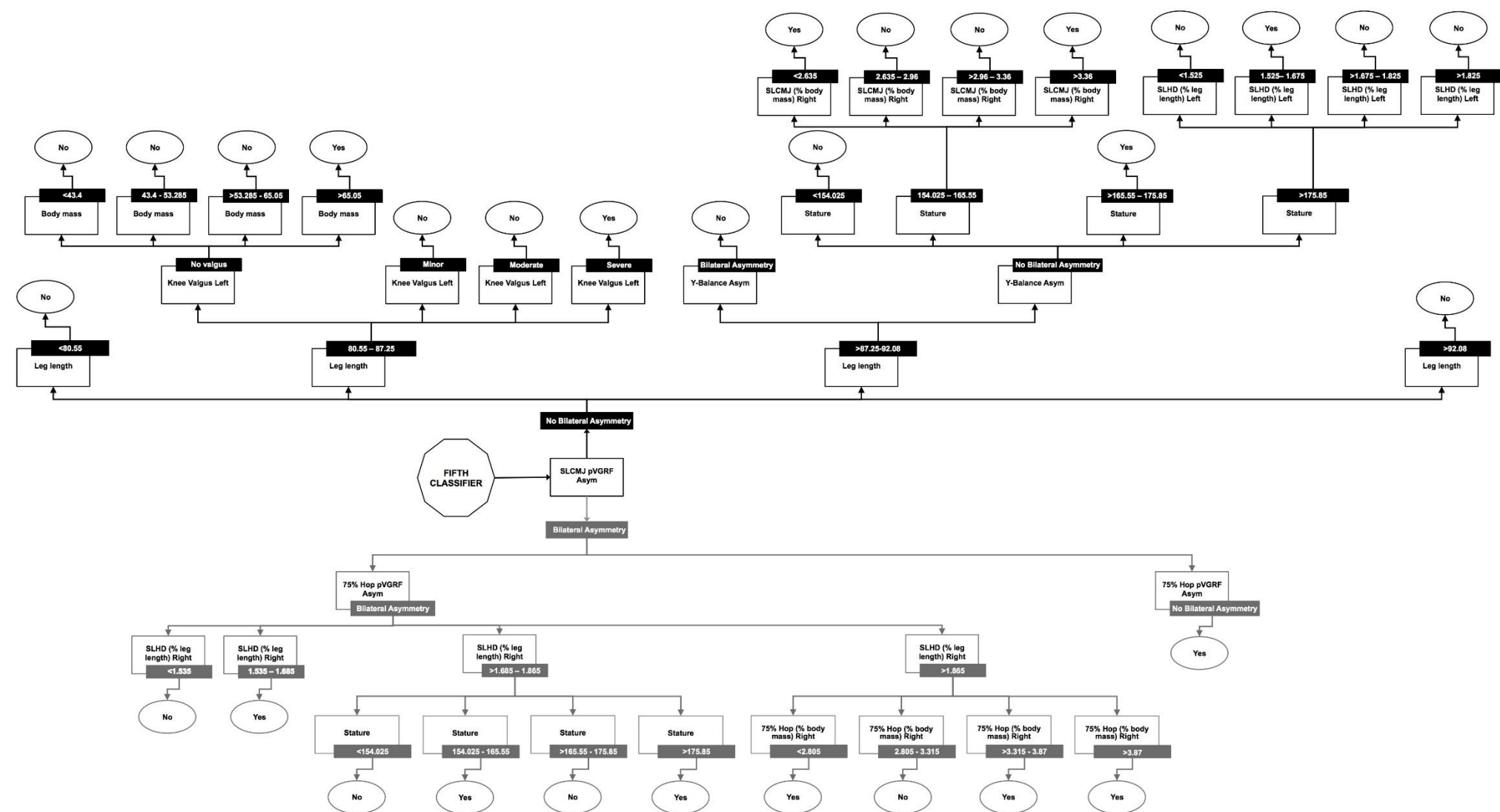
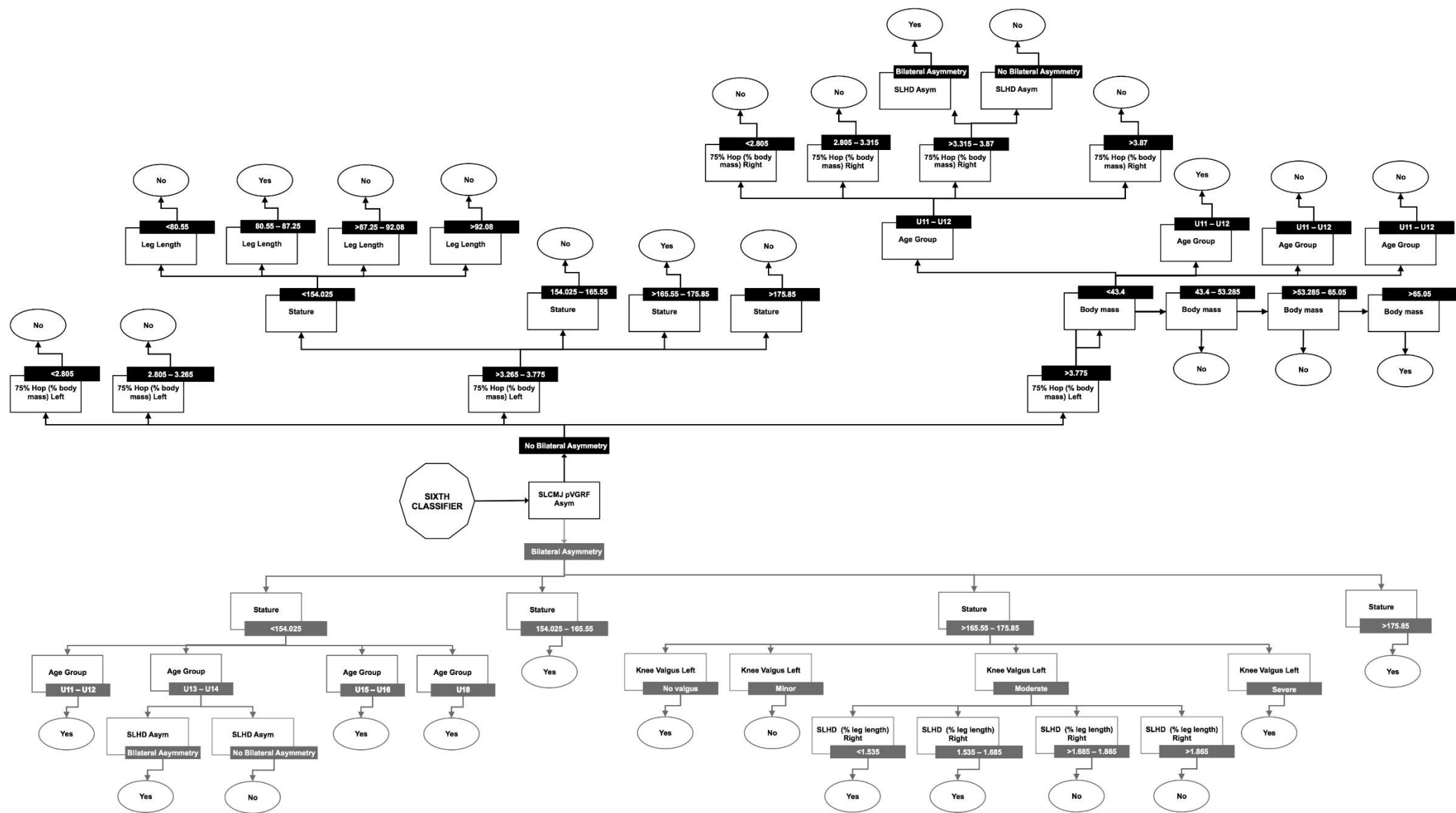


Figure S 6



**Figure S 7**

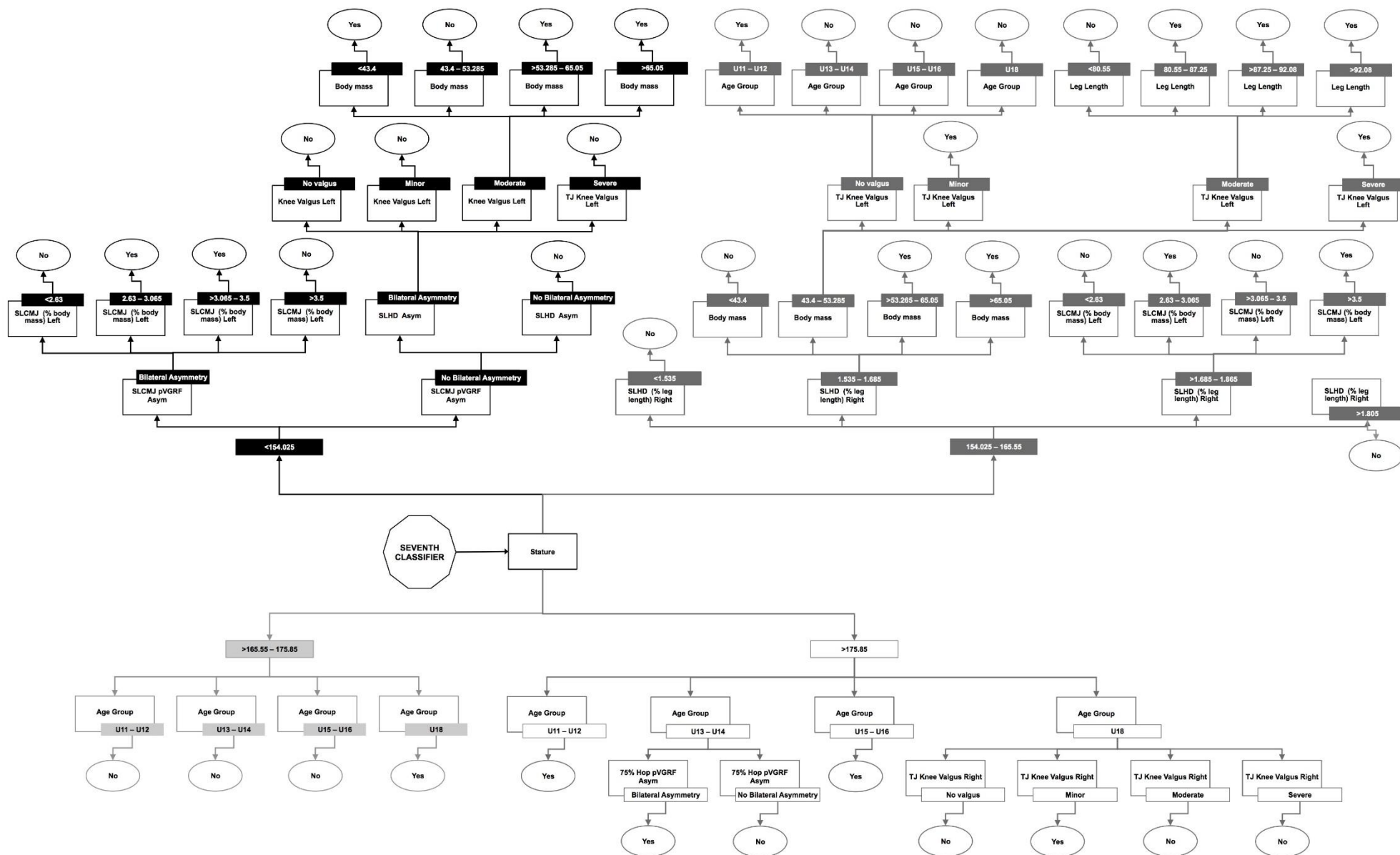
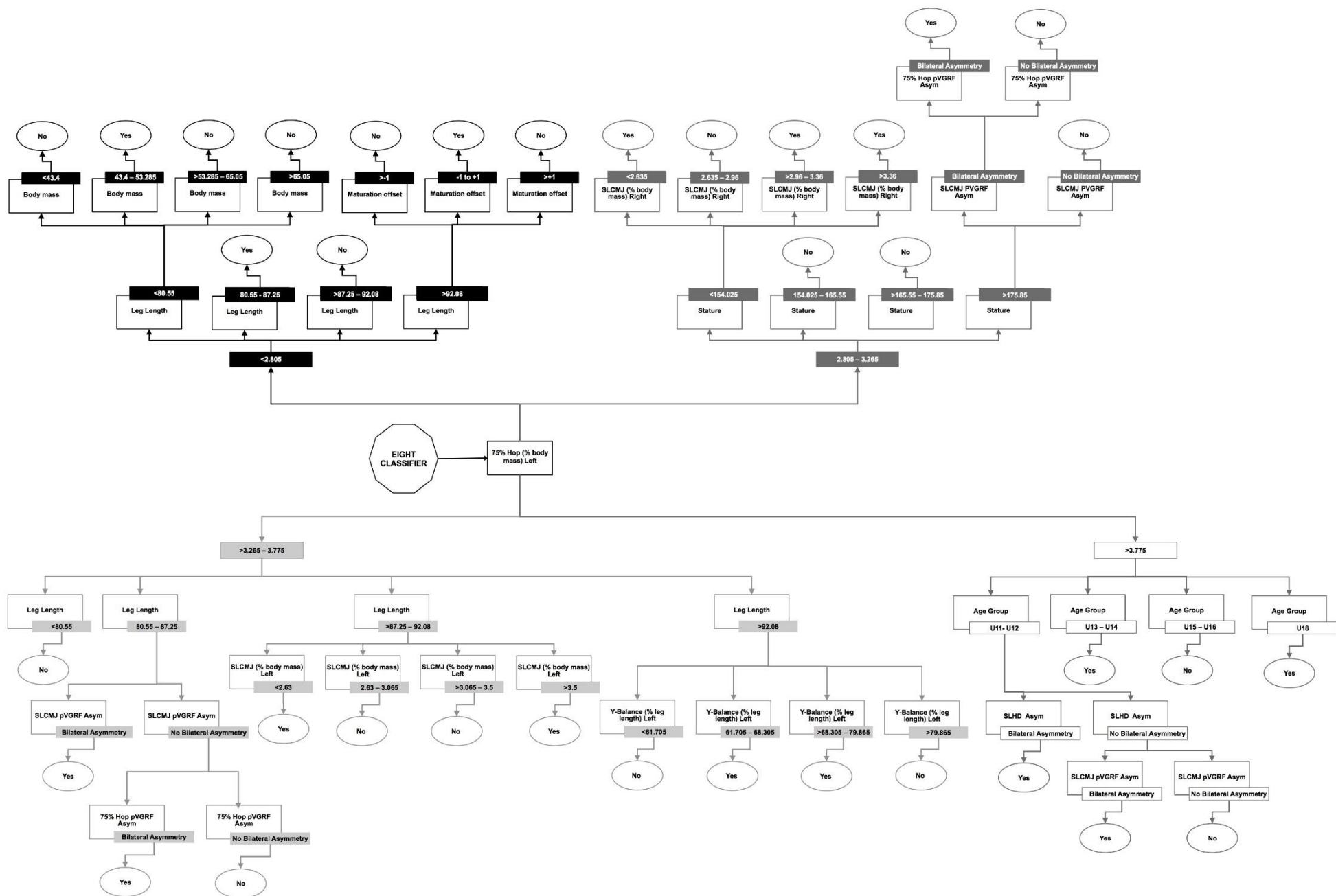


Figure S 8



### Figure S 9



