



Predicting Soccer Injuries with Machine Learning



Jennifer Lawless

DASC 9311: Data Science Project

September 17, 2025



The Problem

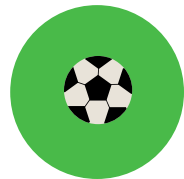
	Core Problem	<ul style="list-style-type: none">• Predict whether a soccer player will suffer an injury in their next season
	Why This Matters	<ul style="list-style-type: none">• Player safety and career longevity• Identifying key injury risk factors• Preventative healthcare• Training program optimization
	Technical Challenges	<ul style="list-style-type: none">• Dataset Size: Limited samples increase overfitting risk• Generalization: Model must work on unseen players• Feature Engineering/Selection: Must avoid multicollinearity and data leakage
	Goals	<ul style="list-style-type: none">• Primary Goal: >85% accuracy with strong generalization• Secondary Goal: Outperform existing Kaggle solutions• Practical Goal: Provide actionable insights for injury prevention

Dataset Overview



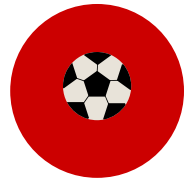
Samples

800 Chinese university soccer players aged 18-24 participating in collegiate and provincial leagues



Features

18 features including physical characteristics, soccer-specific metrics, physical fitness assessments, lifestyle factors, and training compliance



Target

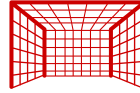
Injury_Next_Season: binary classification value where an injury is defined as training/competition-related injury causing ≥ 7 consecutive days of absence



Data Quality

Multi-source verification (medical records, coach reports, student surveys), **well-balanced** dataset, and **minimal missing data**

Solution Strategy



Comprehensive EDA	Feature Engineering	Model Development	Validation Framework	Multi-Metric Evaluation
<ul style="list-style-type: none">Perform an analysis of injury patterns by position, fitness level, and lifestyle factors	<ul style="list-style-type: none">Create position-relative metrics	<ul style="list-style-type: none">Start with simpler baseline models and then progress to more sophisticated ensembles	<ul style="list-style-type: none">Use nested cross-validation to prevent overfitting	<ul style="list-style-type: none">In addition to Accuracy Rate, use Recall, Precision, F1-Score, and AUC-ROC

Modeling Approach



1

Baseline Models

Logistic Regression, Random Forest XGBoost

2

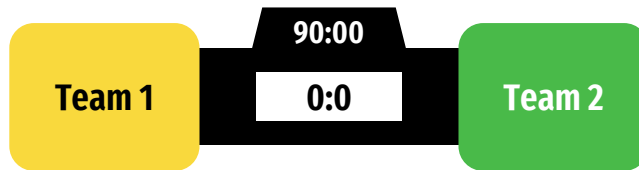
Advanced Models

TabNet, Gaussian Process Classifier, SVM, LightGBM

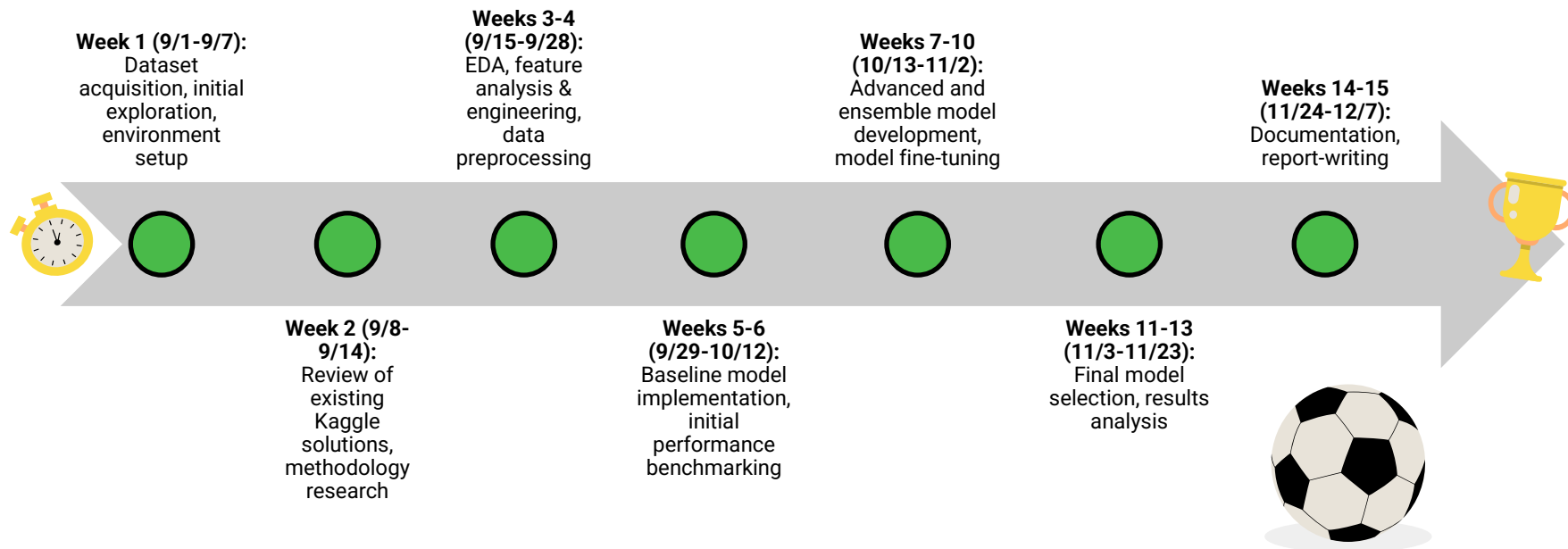
3

Ensemble Methods

Stacking Ensemble, Voting Classifier, Bayesian Model Averaging



Project Timeline



Accomplishments to Date (through 9/17)



Kaggle Notebook Review

Analyzed 8+ existing Kaggle notebooks for this dataset to see that XGBoost achieves 85%+ accuracy in almost all solutions



Environment Setup

Configured Python with necessary packages (scikit-learn, XGBoost, PyTorch, TabNet)



Planning

Designed my model portfolio, established timelines and milestones



Exploratory Data Analysis

Starting basic EDA on the dataset this week



Thank You!

