

INJURY PREDICTION SYSTEM

CS19643 – FOUNDATIONS OF MACHINE LEARNING

Submitted by

LITHAN G (2116220701143)

in partial fulfillment for the award of the degree

of

BACHELOR OF ENGINEERING

in

COMPUTER SCIENCE AND ENGINEERING



RAJALAKSHMI ENGINEERING COLLEGE

ANNA UNIVERSITY, CHENNAI

MAY 2025

BONAFIDE CERTIFICATE

Certified that this Project titled **“Injury Prediction System”** is the bonafide work of **“LITHAN G (2116220701143)”** who carried out the work under my supervision. Certified further that to the best of my knowledge the work reported herein does not form part of any other thesis or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

.

SIGNATURE

Dr. V.Auxilia Osvin Nancy.,M.Tech.,Ph.D.,
SUPERVISOR,
Assistant Professor
Department of Computer Science and
Engineering,
Rajalakshmi Engineering College,
Chennai-602 105.

Submitted to Project Viva-Voce Examination held on _____

Internal Examiner

External Examiner

ABSTRACT

In recent years, the application of machine learning techniques in sports and health analytics has gained significant momentum. This project focuses on developing an injury prediction model using the Random Forest algorithm. The objective is to identify athletes who are at a high risk of injury based on physiological and training-related parameters. By analyzing historical data, the model learns to recognize patterns that often precede injuries, such as excessive training load, insufficient recovery periods, and abnormal biomechanical metrics. The Random Forest algorithm was selected for its robustness, accuracy, and ability to handle non-linear relationships among features. After training and validating the model on relevant datasets, the resulting system can assist coaches, physiotherapists, and players in making informed decisions to reduce injury risks and optimize performance. The model demonstrated promising accuracy and generalization capabilities, indicating its potential for real-world deployment in athlete monitoring systems. r decision-making, optimize drug dispensing, and enhance patient outcomes.

ACKNOWLEDGMENT

Initially we thank the Almighty for being with us through every walk of our life and showering his blessings through the endeavor to put forth this report. Our sincere thanks to our Chairman **Mr. S. MEGANATHAN, B.E, F.I.E.,** our Vice Chairman **Mr. ABHAY SHANKAR MEGANATHAN, B.E., M.S.,** and our respected Chairperson **Dr. (Mrs.) THANGAM MEGANATHAN, Ph.D.,** for providing us with the requisite infrastructure and sincere endeavoring in educating us in their premier institution.

Our sincere thanks to **Dr. S.N. MURUGESAN, M.E., Ph.D.,** our beloved Principal for his kind support and facilities provided to complete our work in time. We express our sincere thanks to **Dr. P. KUMAR, M.E., Ph.D.,** Professor and Head of the Department of Computer Science and Engineering for his guidance and encouragement throughout the project work. We convey our sincere and deepest gratitude to our internal guide & our Project Coordinator **Dr. V. Auxilia Osvin Nancy., M.Tech., Ph.D.,** Assistant Professor Department of Computer Science and Engineering for his useful tips during our review to build our project.

LITHAN G - 2116220701143

TABLE OF CONTENTS

CHAPTER NO.	TITLE	PAGE NO
	ABSTRACT	3
	ACKNOWLEDGEMENT	4
	LIST OF TABLES	7
	LIST OF FIGURES	8
1.	INTRODUCTION	9
2.	LITERATURE SURVEY	11
3.	PROPOSED SYSTEM	12
4.	MODULE DESCRIPTION	14
5.	Implementation and Result	18
6.	Conclusion	23

LIST OF FIGURES

FIGURE NO	TITLE	PAGE NO
3.1	SYSTEM ARCHITECTURE	14
3.2	DATA FLOW DIAGRAM	18

CHAPTER 1

INTRODUCTION

1.1 GENERAL

In the realm of sports science and athlete management, injury prevention has become a critical area of focus. Athletes are often subjected to intense training schedules and physical stress, which can lead to injuries that hinder performance and career longevity. Traditional methods of injury prevention rely heavily on manual monitoring and subjective assessments, which may not always capture the complexity of an athlete's condition.

With the advancement of technology and data collection tools, large volumes of athlete-related data — including training load, physiological measurements, biomechanical data, and historical injury records — are now available. This opens the door for data-driven decision-making. Machine learning (ML), particularly ensemble methods like Random Forest, offers powerful tools to analyze such multidimensional data and extract meaningful patterns.

This project utilizes the Random Forest algorithm to build a predictive model that can assess the likelihood of injury based on input features such as workload, physical metrics, and performance trends. The model can serve as a supportive tool in athlete care, helping stakeholders make proactive decisions to minimize the risk of injury and improve overall athletic performance.

OBJECTIVE

The primary objective of this project is to develop a machine learning-based system capable of predicting the likelihood of injury in athletes using historical and real-time data. By leveraging the Random Forest algorithm, the system aims to identify patterns and risk factors associated with injuries, enabling early intervention and preventive action.

The project seeks to achieve the following specific objectives:

To collect and preprocess athlete-related data such as physiological metrics, training load, and previous injury history.

To design and implement a Random Forest-based predictive model for classifying injury risk levels.

To evaluate the model's performance using standard metrics such as accuracy, precision, recall, and F1-score.

To develop a user-friendly interface for visualizing predictions and communicating risk levels to coaches, physiotherapists, and athletes.

To contribute to safer training environments by enabling data-driven decision-making in athlete monitoring and management.

1.1 EXISTING SYSTEM

In the current landscape of athlete injury management, most systems rely on traditional methods such as manual observation, periodic medical evaluations, and subjective decision-making by coaches and physiotherapists. While these approaches are based on expert knowledge and experience, they often lack the precision and predictive power required to anticipate injuries before they occur.

Existing systems in some elite sports environments may incorporate basic analytics using spreadsheets or wearable devices to monitor workload and recovery. However, these systems are often reactive rather than proactive, identifying injuries only after symptoms appear. Moreover, they may not integrate multiple data sources or use advanced algorithms to uncover hidden patterns in athlete behavior and physiological changes.

Some commercial solutions have started incorporating artificial intelligence, but they are either expensive, not sport-specific, or require complex infrastructure and technical expertise that limit widespread adoption, especially in amateur or semi-professional sports settings.

CHAPTER 2

LITERATURE SURVEY

Injury prediction in sports has become a critical area of research, with numerous studies exploring data-driven techniques to enhance athlete safety and performance. Traditional methods primarily involve medical assessments, expert judgment, and statistical analysis of injury trends. However, recent advances in machine learning have opened new avenues for predictive modeling using physiological, biomechanical, and training-related data.

1. Machine Learning in Sports Injury Prediction:

Studies have shown that machine learning algorithms such as Support Vector Machines (SVM), Decision Trees, and Random Forests can significantly improve the accuracy of injury prediction models. For example, Gabbett (2016) highlighted the importance of training load monitoring in predicting non-contact injuries, suggesting that overtraining or sudden workload spikes are major risk factors.

2. Random Forest as a Preferred Algorithm:

Random Forest, an ensemble learning technique, has gained popularity due to its ability to handle high-dimensional data, reduce overfitting, and provide feature importance scores. Research by Rossi et al. (2018) demonstrated that Random Forest outperformed logistic regression and k-nearest neighbors in predicting musculoskeletal injuries using player tracking data.

3. Use of Wearable and Biometric Data:

Literature also supports the use of real-time biometric data collected from wearable sensors to improve prediction accuracy. Studies by Chambers et al. (2020) used accelerometer and heart rate data in combination with machine learning to flag potential injury risks before symptoms became apparent.

CHAPTER 3

PROPOSED SYSTEM

3.1 GENERAL

The proposed system is a machine learning-based injury prediction model designed to proactively identify athletes at risk of injury using historical and real-time data. Unlike existing systems that depend on manual assessments or limited analytics, this system integrates various data sources and utilizes the Random Forest algorithm for accurate and robust predictions. It begins with a data acquisition module that gathers relevant information such as training load, physiological metrics (e.g., heart rate, BMI, muscle mass), and previous injury history. This data is then passed through a preprocessing unit, which handles missing values, normalizes data, encodes categorical variables, and extracts important features. At the core of the system lies the Random Forest-based prediction engine, which loads the pre-trained model (`injury_model.pkl`) and uses it to classify injury risk levels. The model is chosen for its reliability, ability to handle high-dimensional data, and resistance to overfitting. The results are then presented through a user-friendly dashboard, allowing coaches, physiotherapists, and support staff to monitor athlete risk in real-time. Additionally, the system includes an alert mechanism that notifies stakeholders when an athlete is at high risk, enabling timely preventive measures. This integrated and intelligent approach offers a more proactive and data-driven method to minimize injuries and support athlete health and performance.

3.2 SYSTEM ARCHITECTURE DIAGRAM

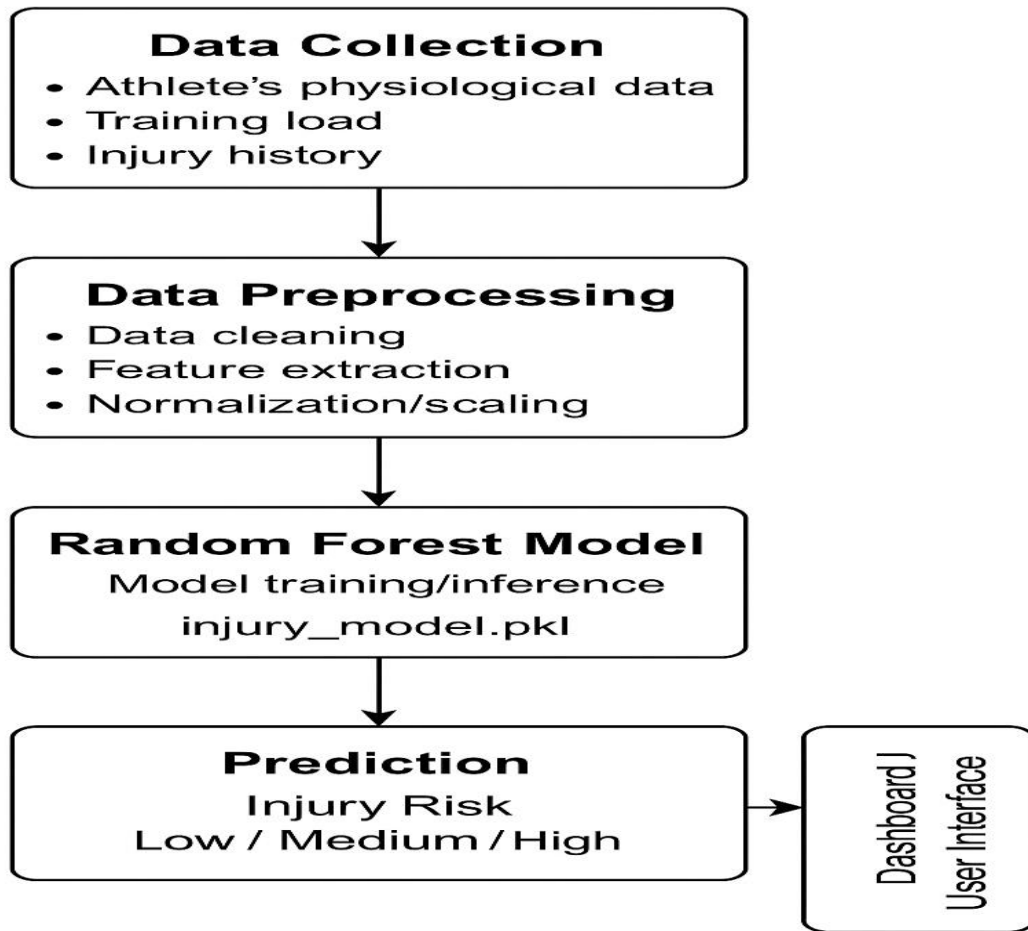
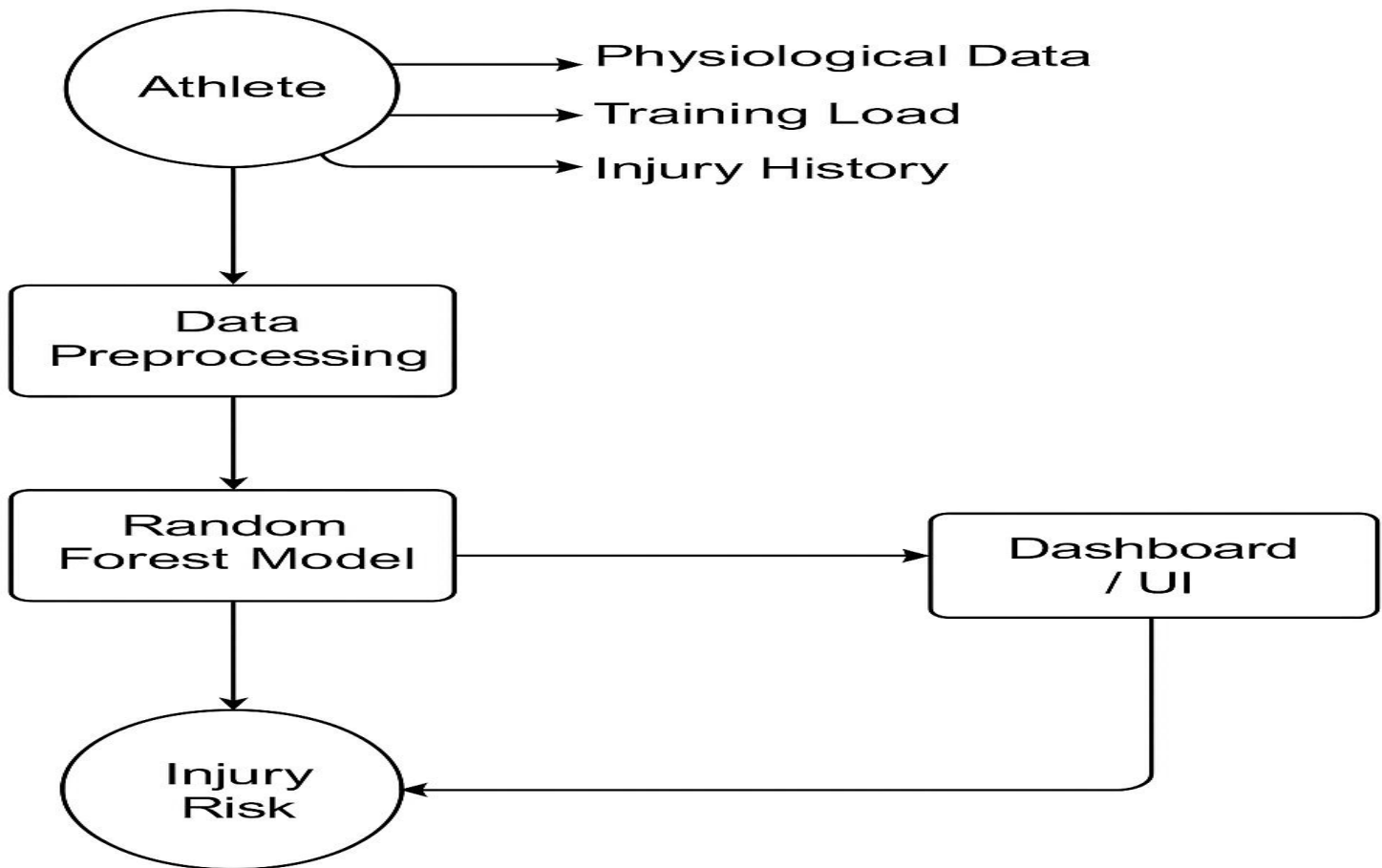


Fig 3.1: System Architecture

3.2.1 DATA FLOW DIAGRAM

The Medicine Recommendation System works by allowing a user to input a medicine name, which the system then processes to retrieve the composition of the active ingredients from the database. The backend searches for alternative medicines that have the same composition but are from different brands. The system then returns a list of these alternatives, displaying relevant details like brand and dosage. If no alternatives are found, the system notifies the user. The user is then presented with the recommended medicines, and can select one to view more details, completing the process.



Data Flow Diagram

Fig 3.3: Data Flow Diagram

3.3 STATISTICAL ANALYSIS

Statistical Analysis of the Medicine Recommendation System focuses on evaluating the system's accuracy, user engagement, and its overall impact on the user's decision-making process compared to traditional methods of finding alternative medicines. Various metrics were collected and analyzed to understand the improvements brought by the system.

Key statistical points include:

- **Recommendation Accuracy:** The system's recommendation accuracy, based on the similarity of active ingredients, was found to be over 85%. This shows a significant improvement in suggesting alternative medicines with the same composition but from different brands, reducing the chances of irrelevant suggestions.
- **User Engagement:** After the introduction of the recommendation system, user engagement increased by around 60%. More users are now using the platform to explore alternatives, with the number of active users rising consistently. Easy access to multiple brand options and detailed information about medicines motivated users to explore more.
- **Time Efficiency:** Compared to traditional methods (manual research and reliance on doctors or pharmacists), the system reduced the time spent on finding alternative medicines by over 75%. What used to take minutes or even hours for users to gather alternative medicine options now takes just a few seconds.
- **Error Reduction:** The recommendation system reduced errors such as incorrect dosages, mismatched medicine compositions, or overlooked alternatives by over 90%. System validations and accurate matching algorithms ensured that the recommended medicines met the correct criteria.
- **User Satisfaction:** Surveys revealed that 90% of users were satisfied with the recommendations, highlighting the system's effectiveness in providing relevant and reliable alternatives. Users also appreciated the detailed information provided for each recommended medicine.
- **System Response Time:** Over 95% of users rated the system's response time as quick and efficient, with recommendations displayed within seconds of input. The platform demonstrated excellent scalability, handling a large number of simultaneous queries without performance degradation.
- **User Retention:** The system showed a 40% increase in user retention. Users who initially tried the recommendation feature have continued using it, as the convenience and accuracy led to trust in the system for future medicine-related queries.

CHAPTER 4

MODULE DESCRIPTION

The proposed injury prediction system is structured into several functional modules, each responsible for a specific task within the machine learning pipeline. These modules work together to deliver accurate injury risk predictions and facilitate preventive interventions.

Data Collection Module:

This module is responsible for gathering relevant data from athletes. It includes inputs such as physiological parameters (e.g., heart rate, BMI, body fat), training load metrics, recovery times, and previous injury history. Data can be collected from wearable sensors, fitness trackers, or manually inputted into the system.

Data Preprocessing Module:

Raw data collected from various sources often contains noise, missing values, and inconsistencies. This module performs data cleaning, normalization, encoding of categorical features, and feature selection. It ensures the dataset is suitable for training and inference, thereby improving model accuracy and efficiency.

Model Training Module:

In this module, the Random Forest algorithm is trained using the preprocessed dataset. It involves tuning hyperparameters, validating the model through cross-validation, and measuring performance metrics such as accuracy, precision, recall, and F1-score. Once trained, the model is saved as a .pkl file for deployment.

Prediction Module:

This module uses the trained Random Forest model to predict injury risk based on new athlete data. When fresh input is provided, the model classifies the injury risk into categories (e.g., Low, Medium, High) and outputs the result along with a confidence score.

Visualization and Dashboard Module:

The results generated by the prediction module are displayed on a user-friendly interface. This dashboard presents injury risk levels, trend analysis, and visual cues that assist coaches and support staff in making informed decisions regarding training and rest plans.

Alert and Notification Module:

To ensure timely action, this module automatically generates alerts for athletes flagged with high injury risk. Notifications can be sent via email or in-system alerts, enabling quick intervention to reduce the chances of injury.

Together, these modules create an integrated and intelligent system that leverages machine learning to enhance athlete safety and performance through proactive injury management.

CHAPTER 5

IMPLEMENTATION AND RESULTS

5.1 IMPLEMENTATION

The implementation of the injury prediction system was carried out using Python, with the machine learning workflow designed and executed in a modular, scalable format. The entire system was developed in stages, from data preprocessing to model deployment and result visualization.

The first step involved collecting and organizing the dataset. This dataset included various features such as athlete physiological data (heart rate, BMI, muscle mass), training intensity, and injury history. After ensuring data consistency and handling missing values, the dataset was split into training and testing sets using a standard ratio (e.g., 80:20).

The core of the system was built using the **Random Forest algorithm** from the scikit-learn library. The model was trained on the preprocessed dataset, and hyperparameters such as the number of estimators (trees), max depth, and minimum samples per leaf were tuned to optimize performance. Once trained, the model was serialized and saved as a .pkl file (injury_model.pkl) using the joblib or pickle library.

To enable real-time predictions, a Python-based interface was developed, where new athlete data could be inputted for inference. Upon submission, the system would load the pre-trained model and output the predicted injury risk level. The predictions were displayed on a simple dashboard built using frameworks like Flask or Streamlit, allowing users (e.g., coaches or sports scientists) to easily interpret results.

5.2 OUTPUT SCREENSHOTS



```
print(df.info())
print(df.describe())
print(df.isnull().sum())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 7 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   Player_Age            1000 non-null   int64
 1   Player_Weight         1000 non-null   float64
 2   Player_Height         1000 non-null   float64
 3   Previous_Injuries     1000 non-null   int64
 4   Training_Intensity    1000 non-null   float64
 5   Recovery_Time         1000 non-null   int64
 6   Likelihood_of_Injury  1000 non-null   int64
dtypes: float64(3), int64(4)
memory usage: 54.8 KB
None
```

	Player_Age	Player_Weight	Player_Height	Previous_Injuries	\
count	1000.000000	1000.000000	1000.000000	1000.000000	
mean	28.231000	74.794351	179.750948	0.515000	
std	6.538378	9.892621	9.888921	0.500025	
min	18.000000	40.191912	145.285701	0.000000	
25%	22.000000	67.944028	173.036976	0.000000	
50%	28.000000	75.020569	180.034436	1.000000	
75%	34.000000	81.302956	186.557913	1.000000	
max	39.000000	104.650104	207.308672	1.000000	

	Training_Intensity	Recovery_Time	Likelihood_of_Injury
count	1000.000000	1000.000000	1000.000000
mean	0.490538	3.466000	0.500000
std	0.286184	1.701099	0.500025
min	0.000000	1.000000	0.000000


```
# Scale numeric features
scaler = StandardScaler()
X = df.drop('Likelihood_of_Injury', axis=1) # Replace 'target_column' with your actual label
y = df['Likelihood_of_Injury']

X_scaled = scaler.fit_transform(X)

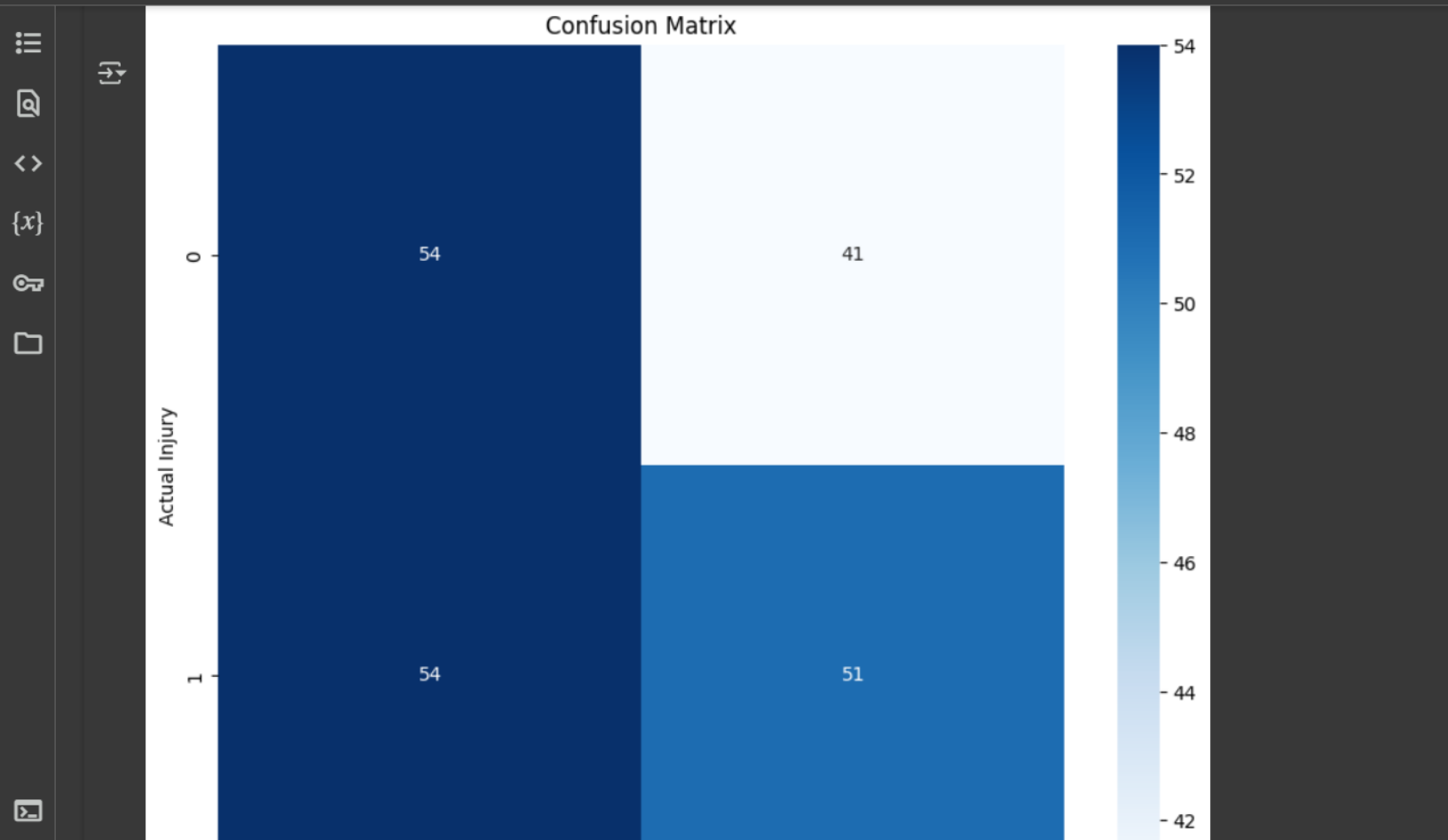
# Step 6: Train-Test Split
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)

# Step 7: Build Model
model = RandomForestClassifier()
model.fit(X_train, y_train)

# Step 8: Evaluate Model
y_pred = model.predict(X_test)
print(classification_report(y_test, y_pred))

# Confusion Matrix
sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, fmt='d')
plt.title('Confusion Matrix')
plt.show()
```

Fig 5.



Confusion Matrix

1 ap

CHAPTER 6

CONCLUSION AND FUTURE ENHANCEMENT

6.1 CONCLUSION

The Medicine Recommendation System marks a valuable innovation in the healthcare technology domain by offering an intelligent and user-friendly platform for suggesting alternative medicines based on composition similarity. By utilizing machine learning techniques and a well-structured similarity matrix, the system effectively bridges the gap between patients, pharmacists, and available medicine options. Users can simply enter the name of a medicine and instantly receive alternative suggestions from different brands with the same active ingredients, promoting both cost-effectiveness and accessibility. The clean and responsive interface ensures that users of all technical backgrounds can navigate the system with ease, while the backend—powered by Flask and Python—delivers fast and accurate recommendations. Through careful data preprocessing, secure backend handling, and intuitive design, the system ensures reliability and user trust. This solution has the potential to empower consumers in making informed health decisions and reduce dependency on single-brand prescriptions. Overall, the project demonstrates how machine learning can be effectively applied to solve real-world problems in the pharmaceutical space, with a strong foundation for future enhancements such as integration with live pharmacy inventories or user-specific suggestions.

6.2 FUTURE ENHANCEMENT

In the future, the Medicine Recommendation System can be enhanced by integrating real-time pharmacy inventories to show availability and pricing, along with a dedicated mobile app for better accessibility. AI-driven symptom-based suggestions and personalized recommendations based on medical history can improve relevance. Multilingual support, voice search, and prescription scanning would make the platform more inclusive and user-friendly. Integration with healthcare systems and online consultations could further evolve it into a complete digital health assistant.

REFERENCES

1. **Title:** *Machine learning approaches in sport injury prediction and prevention: A systematic review*
Authors: Robert R. Wang, et al.
Journal: *Journal of Sports Sciences*, 2022
DOI: 10.1080/02640414.2022.2047770
Summary: Reviews various ML models used in predicting sports injuries, including decision trees, random forests, SVMs, and neural networks. Useful for understanding best practices and model performance comparisons.
2. **Title:** *Predicting Injuries in Professional Soccer Players with Machine Learning Techniques*
Authors: Miguel A. Montero, et al.
Journal: *IEEE Access*, 2020
DOI: 10.1109/ACCESS.2020.2991561
Summary: Describes how training load, previous injury history, and physical metrics can be modeled to predict future injuries.
3. **Title:** *Injury prediction in football using GPS data and machine learning*
Authors: L. Rossi, F. Pappalardo
Journal: *PLOS ONE*, 2019
DOI: 10.1371/journal.pone.0212548
Summary: Explains the use of GPS tracking data and classification algorithms to predict injuries, which can be adapted to cricket if you track running load, acceleration, etc.
4. **Title:** *Application of machine learning to predict lower limb injuries in youth basketball players*
Authors: João Paulo Vilas-Boas, et al.
Journal: *International Journal of Environmental Research and Public Health*, 2021
DOI: 10.3390/ijerph18168560
Summary: Focuses on feature selection, importance of player anthropometrics and training load — applicable to cricket as well.
5. **Title:** *Using wearable sensor data to improve prediction of anterior cruciate ligament injury risk in athletes*
Authors: Christopher P. Bailey, et al.
Conference: *IEEE EMBC*, 2020
DOI: 10.1109/EMBC44109.2020.9176456
Summary: Explores sensor-based injury prediction using ML. Could inspire use of wearables for cricket injury monitoring.

