

# **EXPLAINABLE ARTIFICIAL INTELLIGENCE**

MINI PROJECT REPORT

Submitted By

**ABHINAV R(MCC21CS004)**

**ABHIRAGH A R(MCC21CS005)**

**ALFIN MUHAMMED N S(MCC21CS012)**

**ASWATHY M L(MCC21CS021)**

**to**

The APJ Abdul Kalam Technological University

In partial fulfilment of the requirements for the award of the Degree

of

Bachelor of Technology

In

*Computer Science and Engineering*



**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**MUSALIAR COLLEGE OF ENGINEERING**

**KADAKAM P.O CHIRAYINKEEZHU, TRIVANDRUM 695304**

May 2024

## **DECLARATION**

We undersigned hereby declare that the project report “**EXPLAINABLE ARTIFICIAL INTELLIGENCE**”, submitted for partial fulfilment of the requirements for the award of degree of Master of Technology of the APJ Abdul Kalam Technological University, Kerala is a bonafide work done by us under supervision of Mrs. Labeeba Vahid. This submission represents our ideas in our own words and where ideas or words of others have been included, we have adequately and accurately cited and referenced the original sources. we also declare that we have adhered to ethics of academic honesty and integrity and have not misrepresented or fabricated any data or idea or fact or source in my submission. We understand that any violation of the above will be a cause for disciplinary action by the institute and/or the University and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been obtained. This report has not been previously formed the basis for the award of any degree, diploma or similar title of any other University.

**ABHINAV R(MCC21CS004)**

**ABHIRAGH A R (MCC21CS005)**

**ALFIN MUHAMMED N S (MCC21CS012)**

**ASWATHY M L (MCC21CS021)**

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**  
**MUSALIAR COLLEGE OF ENGINEERING KADAKAM P.O**  
**CHIRAYINKEEZHU, TRIVANDRUM 695304**



**CERTIFICATE**

This is to certify that, the project report entitled '**EXPLAINABLE ARTIFICIAL INTELLIGENCE**', is a Bonafide record of the CSD334 Project presented by **ABHINAV R(MCC21CS004), ABHIRAGH A R(MCC21CS005),ALFIN MUHAMMED N S(MCC21CS012), ASWATHY M L(MCC21CS021)** of Sixth Semester B. Tech. Computer Science & Engineering students, under our guidance and super vision, in partial fulfilment of the requirements for the award of the degree, B. Tech. Computer Science & Engineering of APJ Abdul Kalam Technological University.

Project Guide

Project Coordinator

HOD

**Mrs. Labeeba Vahid**

**Mrs.Aswathy Gandhi**

**Mrs.Remya R S**

Assistant Professor

Assistant Professor

Assistant Professor

## ACKNOWLEDGEMENT

This work would not have been possible without the support of many people. First and the foremost, we give thanks to Almighty God who gave us the inner strength, resource and ability to complete our project successfully.

We would like to thank **Dr. K K Abdul Rasheed**, our Principal, who has provided with the best facilities and atmosphere for the project completion and presentation. We would also like to thank our HOD **Mrs. Remya R S** (Assistant Professor, Department Of Computer Science Engineering) and our project coordinator **Mrs. Aswathy Gandhi**(Assistant Professor, Department Of Computer Science Engineering), our project guide **Mrs. Labeeba Vahid** (Assistant Professor, Department Of Computer Science Engineering) and all of our faculties for the help extended and also for the encouragement and support given to us while doing the project.

We would like to thank my dear friends for extending their cooperation and encouragement throughout the project work, without which we would never have completed the project this well.

Thank you all for your love and also for being very understanding.

# CONTENTS

| Contents                                         | Page No |
|--------------------------------------------------|---------|
| LIST OF FIGURES                                  | i       |
| LIST OF TABLES                                   | ii      |
| LIST OF ABBREVIATIONS                            | iii     |
| ABSTRACT                                         | iv      |
| CHAPTER 1 – INTRODUCTION                         | 1       |
| 1.1 The Enigma Of Black Box Models               |         |
| 1.2 Demystifying The Machine: The Rise Of Xai    |         |
| 1.3 Objectives: An Exploration Of Xai Techniques |         |
| 1.4 Expected Benefits                            |         |
| 1.5 Project Scope                                |         |
| CHAPTER 2 - LITERATURE SURVEY                    | 5       |
| 2.1 Interpretable Matrix Factorization           |         |
| 2.2 Explainable Systems: Survey And Challenges   |         |
| 2.3 Towards User-Centric Explainable Systems     |         |
| 2.4 Interpretable Neural Collaborative Filtering |         |
| 2.5 Human-Centered Explainable Ai: A Survey      |         |
| CHAPTER 3 - EXISTING SYSTEM                      | 9       |
| 3.1 The Realm Of Non-Explainable Models          |         |
| 3.2 A Shift Towards Transparency                 |         |

|                                                 |    |
|-------------------------------------------------|----|
| CHAPTER 4 - PROPOSED SYSTEM                     | 10 |
| 4.1 Choosing The Domain: A Practical Scenario   |    |
| 4.2 Tailoring XAI For The Model                 |    |
| 4.3 Building The Integrated System              |    |
| CHAPTER 5 - MODULE DESCRIPTION                  | 11 |
| 5.1 Module 1: Data Collection And Preprocessing |    |
| 5.2 Module 2: Model Building And Training       |    |
| 5.3 Module 3: Explainable Ai Integration        |    |
| 5.4 Module 4: Evaluation And Refinement         |    |
| CHAPTER 6 – REQUIREMENTS                        | 13 |
| 6.1 Software Requirements                       |    |
| 6.2 Hardware Requirements                       |    |
| CHAPTER 7 – SYSTEM DESIGN                       | 14 |
| CHAPTER 8 - SAMPLE CODE                         | 15 |
| CHAPTER 9 - RESULT AND OUTPUT                   | 22 |
| CHAPTER 10 - FUTURE ENHANCEMENT                 | 27 |
| 10.1 Expanding Explanations                     |    |
| 10.2 Beyond Real Estate And Titanic Data        |    |
| 10.3 Human-AI Collaboration                     |    |
| CHAPTER 11 – CONCLUSION                         | 28 |
| CHAPTER 12 – REFERENCES                         | 29 |

## **LIST OF FIGURES**

| <b>FIG NO</b> | <b>TITLE</b>                      | <b>Page no</b> |
|---------------|-----------------------------------|----------------|
| 7.1           | System Design Overview            | 14             |
| 9.1           | Home Page                         | 22             |
| 9.2           | House Price Prediction Model      | 23             |
| 9.3           | Titanic Survivor Prediction Model | 24             |
| 9.4           | SHAP Explanation                  | 25             |
| 9.5           | LIME Explanation                  | 26             |

## **LIST OF TABLES**

| <b>FIG NO</b> | <b>TITLE</b>      | <b>Page no</b> |
|---------------|-------------------|----------------|
| 2.1           | Literature Survey | 7              |



## **LIST OF ABBREVIATIONS**

AI – Artificial Intelligence

XAI – Explainable Artificial Intelligence

CNN - Convolutional Neural Network

RNN - Recurrent Neural Network

SHAP - SHapley Additive exPlanations

LIME - Local Interpretable Model-agnostic Explanations

RAM – Random Access Memory

VCS – Version Control System

# ABSTRACT

As Artificial Intelligence (AI) permeates our lives, influencing everything from loan approvals to medical diagnoses, a critical question emerges: can we understand how AI models arrive at their decisions? Enter Explainable AI (XAI), a transformative field that strives to demystify the inner workings of these complex algorithms. This project delves into the practical application of XAI techniques, aiming to bridge the gap between the opaque nature of AI models and human comprehension.

By integrating XAI methodologies, we can illuminate the "why" behind AI predictions. This newfound transparency fosters trust and empowers users with a deeper understanding of the rationale governing AI decisions. Imagine a world where loan applicants can comprehend the factors influencing loan approval/rejection, or where medical professionals gain insights into the reasoning behind AI-powered diagnostic tools. XAI empowers informed decision-making and fosters collaboration between humans and AI.

This project goes beyond theoretical exploration. It equips you with the tools and knowledge to implement XAI in your own machine learning endeavors. We explore practical considerations such as hardware and software requirements, along with valuable resources to guide your journey. Furthermore, the project investigates the potential for XAI across various domains, from real estate prediction to loan approval systems and beyond. By demonstrating the applicability of XAI in diverse scenarios, we pave the way for its widespread adoption.

Ultimately, this project aspires to contribute to a future where AI development is guided by principles of responsibility and transparency. By fostering a future where humans and AI work in concert, we can harness the immense potential of AI for the betterment of society. This exploration of XAI is not just about unveiling the black box; it's about unlocking a future where AI can be a trusted and powerful force for good.

# CHAPTER 1

## INTRODUCTION

From revolutionizing facial recognition software to personalizing recommendations on streaming platforms, Artificial Intelligence (AI) has undeniably transformed our world. These intelligent systems are constantly evolving, pushing the boundaries of automation and decision-making across diverse fields. However, as AI models become increasingly powerful and complex, a critical challenge emerges: the lack of transparency in their decision-making processes. Many of these models, often referred to as "black boxes," excel at generating accurate predictions. For instance, an AI system designed for loan approval might consistently determine creditworthiness with remarkable accuracy. Yet, the inner workings of this system remain a mystery. We may know that the loan is approved, but we lack a clear understanding of the specific factors that influenced this decision. This lack of transparency raises a multitude of concerns, hindering trust, accountability, and responsible development of AI technologies.

One of the most pressing concerns surrounding black box models is the potential for bias. AI models are trained on vast datasets, and these datasets can unconsciously reflect societal biases present in the real world. If the data used to train a loan approval system disproportionately favors applicants with certain demographics or financial backgrounds, the resulting model might perpetuate those biases in its decision-making. Without understanding how the model arrives at its conclusions, it becomes incredibly difficult to identify and address such biases. This lack of transparency can lead to discriminatory outcomes, hindering equal access to opportunities and resources.

Furthermore, the lack of interpretability in black box models poses challenges for accountability and regulatory oversight. Suppose a self-driving car makes an unexpected decision on the road, potentially leading to an accident. Without understanding the reasoning behind its actions, it becomes difficult to determine the cause of the accident and hold the appropriate party accountable.

## 1.1 THE ENIGMA OF BLACK BOX MODELS

The rise of Artificial Intelligence (AI) has been nothing short of phenomenal. From the intelligent assistants that seamlessly integrate into our daily lives to the sophisticated algorithms powering self-driving cars and medical diagnostics, AI is transforming how we interact with the world and make critical decisions. However, as these models become increasingly complex and powerful, a crucial question arises: how do they actually arrive at their outputs? Many AI models, often referred to as "black boxes," excel at generating accurate predictions. For example, an AI system used for loan approval might consistently determine creditworthiness with remarkable accuracy. Yet, the inner workings of this system remain shrouded in mystery. We know whether a loan is approved or denied, but the specific factors influencing this decision remain hidden from view. This lack of transparency in black box models poses a significant challenge, hindering trust, accountability, and ultimately, the responsible development of AI technologies.

## 1.2 DEMYSTIFYING THE MACHINE: THE RISE OF XAI

The limitations of black box models have spurred the development of a crucial counterpoint: Explainable AI (XAI). XAI techniques aim to illuminate the decision-making processes of AI models, transforming them from opaque oracles into systems with a degree of transparency. By shedding light on the rationale behind model predictions, XAI offers a multitude of benefits that address the challenges outlined in the previous section.

- **Fostering Trust and Transparency:** As discussed earlier, the lack of transparency in black box models hinders trust in AI systems. XAI techniques bridge this gap by providing explanations for model outputs. These explanations can take various forms, such as highlighting the most influential features in a decision, visualizing the model's internal logic, or offering counterfactual examples that demonstrate how changing specific inputs might affect the prediction.
- **Detecting and Mitigating Bias:** The susceptibility of black box models to bias is a critical concern. XAI techniques can be instrumental in identifying and mitigating these biases. By analyzing how different features contribute to model predictions, XAI methods can expose unforeseen biases lurking within the training data or the model architecture itself.

## 1.3 OBJECTIVES: AN EXPLORATION OF XAI TECHNIQUES

This project doesn't just focus on theoretical exploration. A core objective is to integrate the chosen XAI techniques and models into a practical application, demonstrating their value in a real-world scenario. This could involve building a system that leverages XAI to explain the predictions of a machine learning model used in a specific domain. For instance, we might develop an application that utilizes XAI to explain loan approval decisions made by an AI-powered lending system. By integrating XAI into this real-world application, we aim to showcase its practical utility and demonstrate how it can foster transparency and trust in AI-driven decision-making processes.

## 1.4 EXPECTED BENEFITS:

By achieving these objectives, this project anticipates several benefits:

- **Enhanced Trust and Transparency:** By demystifying model behavior, XAI can increase trust in AI systems, leading to wider adoption and acceptance.
- **Bias Detection and Mitigation:** Unveiling the factors influencing model decisions can enable the identification and mitigation of potential biases embedded within the data or model training process.
- **Improved Model Development:** Understanding how models arrive at their outputs can guide developers in optimizing model architectures and enhancing their overall performance.
- **Informed User Decisions:** Users can make more informed decisions by gaining insights into the reasoning behind a model's predictions.

## 1.5 PROJECT SCOPE:

- **Model Selection:** We will select a set of diverse machine learning models encompassing different tasks and architectures (e.g., Random Forest for regression, Convolutional Neural Network for image classification).

- **XAI Techniques:** We will explore a range of XAI techniques, such as SHAP explanations for tree-based models and LIME for more complex models.
- **Dataset Selection:** Datasets will be chosen that are relevant to the chosen models and tasks, ensuring compatibility and meaningful explanations.
- **Real-World Application Scenario:** We will define a specific real-world scenario for model and XAI integration, demonstrating its practical utility.

## **CHAPTER 2**

### **LITERATURE SURVEY**

Recommender systems (RS) play a crucial role in today's information-laden world, helping users navigate vast selections of products, movies, music, and more. However, the opaque nature of traditional RS models often leaves users in the dark about why a particular item is recommended. This lack of transparency can hinder trust and user satisfaction. To address this challenge, research in Explainable Artificial Intelligence (XAI) is paving the way for the development of explainable recommender systems (XAI-RS).

#### **2.1 INTERPRETABLE MATRIX FACTORIZATION**

In their work titled "Interpretable Matrix Factorization for Recommender Systems" (2010), Rendle et al. propose a method for interpretable matrix factorization (IMF) in RS. IMF is a popular technique for recommendation, but it can be challenging to understand the rationale behind its predictions. Rendle et al. address this by incorporating predefined features into the model, allowing for a degree of interpretability. These features could represent characteristics of users (e.g., demographics, interests) or items (e.g., genre, brand).

#### **2.2 EXPLAINABLE SYSTEMS: SURVEY AND CHALLENGES**

Wilk et al. (2022) present a broader survey of explainability in RS titled "Explainable Recommendation Systems: Survey and Challenges." Their work highlights the importance of XAI in RS and explores various approaches to explainability, including model-agnostic and model-specific techniques. Model-agnostic techniques are independent of the specific RS model used, while model-specific techniques are tailored to the inner workings of a particular model.

#### **2.3 TOWARDS USER-CENTRIC EXPLAINABLE SYSTEMS**

Xu et al. (2023) delve deeper into the concept of user-centric explainability in their paper "Towards User-Centric Explainable Recommendation Systems." They emphasize the need for explanations

that resonate with users and address potential biases within the model. While some approaches focus on explaining individual recommendations, Xu et al. highlight the importance of explaining the potential for bias within the system. This transparency can foster trust with users by acknowledging the limitations of the model and how biases in the training data could influence recommendations.

## **2.4 INTERPRETABLE NEURAL COLLABORATIVE FILTERING**

He et al. (2023) explore interpretable neural collaborative filtering (NCF) for XAI-RS in their work titled "Interpretable Neural Collaborative Filtering for Explainable Recommendation Systems." NCF is a powerful technique for recommendations, but it can be challenging to interpret its inner workings due to its complex architecture. He et al. propose methods for incorporating interpretability into the NCF model by introducing additional layers or modifying existing ones to provide explanations for model predictions. However, a limitation of their approach is the focus on specific interaction data, such as user ratings or clicks.

## **2.5 HUMAN-CENTERED EXPLAINABLE AI: A SURVEY**

Singh et al. (2023) delve into the concept of human-centered XAI, emphasizing the importance of tailoring explanations to user needs and preferences. They identify a gap in research regarding specific XAI techniques that resonate best with users. While various explainability methods exist (e.g., feature importance, counterfactual explanations, visualizations), it's crucial to understand how these methods translate to user understanding and satisfaction.

**User Studies:** Conduct user studies to evaluate the effectiveness of different XAI techniques in recommender systems. This can involve observing user interactions with explanations, gathering user feedback on clarity and usefulness, and comparing different explanation formats.

**Tailoring Explanations to User Preferences:** Investigate ways to personalize explanations based on user characteristics like technical expertise, preferred level of detail, and desired format (textual, visual, interactive). This personalization can ensure that explanations are not only informative but also engaging and relevant to individual users.



| <b>Serial No</b> | <b>Title</b>                                               | <b>Author</b>        | <b>Limitations</b>                                                                       | <b>Proposed Improvements</b>                                                                                           |
|------------------|------------------------------------------------------------|----------------------|------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------|
| 1.               | Interpretable Matrix Factorization for Recommender Systems | Steffen Rendle et al | Pre-defined features struggle to capture user preferences and changing tastes over time. | Integrate user feedback and dynamic updates to capture evolving preferences and provide more tailored recommendations. |
| 2.               | Explainable Recommendation Systems: Survey and Challenges  | Bartosz Wilk et al   | Focuses heavily on model centric explanations neglecting user preferences and feedback   | Incorporate user preferences and feedback into the explanation process.                                                |

|    |                                                                                     |                    |                                                                                                  |                                                                                              |
|----|-------------------------------------------------------------------------------------|--------------------|--------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------|
| 3. | Towards User-Centric Explainable Recommendation Systems                             | Lei Xu et al       | Primarily focuses on explaining individual recommendations.                                      | Provide explanations for potential biases within the model, promoting transparency and trust |
| 4. | Interpretable Neural Collaborative Filtering for Explainable Recommendation Systems | Xiangnan He et al  | Focus on specific interaction data: The model relies on user-item interactions (ratings, clicks) | Incorporate additional user data: Explore ways to integrate user reviews                     |
| 5. | Human-Centered Explainable AI: A Survey                                             | Sameer Singh et al | Limited focus on technical details and user centrality.                                          | Focus on specific XAI techniques that align well with user-centrality.                       |

Fig 2.1 Literature Survey

## CHAPTER 3

### EXISTING SYSTEM

#### 3.1 THE REALM OF NON-EXPLAINABLE MODELS:

Non-explainable models, often referred to as black boxes, have revolutionized various fields with their ability to generate accurate predictions. From facial recognition software to spam filters, these models operate with remarkable efficiency. However, their inner workings remain shrouded in mystery.

The Drawbacks of Opacity:

- **Erosion of Trust:** Users have difficulty trusting a system whose decision-making process is a mystery. Without understanding the "why" behind a prediction, users may question its fairness and accuracy.
- **Hidden Biases:** Non-explainable models can unknowingly perpetuate biases present in the training data. Since we can't see how the model arrives at its conclusions, identifying and mitigating these biases becomes a significant challenge.
- **Hindered Improvement:** The inability to understand how a model arrives at its outputs hinders efforts to improve its performance. If we can't pinpoint weaknesses or areas for refinement in the model's internal logic, it's difficult to effectively enhance its capabilities.

#### 3.2 A SHIFT TOWARDS TRANSPARENCY:

By embracing XAI, we move beyond the limitations of non-explainable models. XAI empowers users, fosters trust in AI systems, and paves the way for the development of more responsible and reliable AI technologies. The next chapter will delve into the practical application of XAI techniques, showcasing their value in a real-world scenario.

## CHAPTER 4

### PROPOSED SYSTEM

#### 4.1 CHOOSING THE DOMAIN: A PRACTICAL SCENARIO

To showcase the value of XAI, we will select a specific domain and scenario where an AI model plays a crucial role. Once the domain and scenario are chosen, we will select a suitable machine learning model for the task.

#### 4.2 TAILORING XAI FOR THE MODEL:

Building upon the XAI techniques explored, we will select the most appropriate methods for the chosen machine learning model. Here are some factors to consider during this selection process:

- **Model Architecture:** The XAI technique should be compatible with the architecture of the chosen model (e.g., decision trees, neural networks).
- **Desired Explanation Level:** Do we require global explanations for overall model behavior (e.g., SHAP) or local explanations for specific predictions (e.g., LIME)?
- **Interpretability:** The chosen XAI technique should provide explanations that are interpretable by the intended audience (e.g., domain experts, loan applicants).

#### 4.3 BUILDING THE INTEGRATED SYSTEM:

This section will delve into the technical details of integrating the chosen XAI techniques with the selected machine learning model. This might involve:

- **Training the Model:** The machine learning model will be trained on relevant data specific to the chosen domain and scenario.
- **Incorporating XAI Techniques:** The selected XAI methods will be integrated into the system to enable explanation generation during model prediction.

## CHAPTER 5

### MODULE DESCRIPTION

#### 5.1 MODULE 1: DATA COLLECTION AND PREPROCESSING

➤ **Data Collection:**

- **House Price Prediction:** Gather real estate data containing features like square footage, number of bedrooms, location, and sale prices. Sources could include online real estate listings or public datasets.
- **Titanic Survivor Prediction:** Acquire the Titanic passenger dataset, which includes information like passenger class, age, gender, and survival outcomes. This data is readily available from public repositories.

➤ **Data Preprocessing:**

- Clean and format the data to ensure consistency and address missing values or outliers. This might involve data imputation, normalization, or feature engineering.
- Explore the data to identify potential relationships between features and target variables (house price, survival probability). This initial exploration can provide insights for selecting appropriate XAI techniques later.

#### 5.2 MODULE 2: MODEL BUILDING AND TRAINING

➤ **House Price Prediction:**

- Choose a suitable regression model, such as Random Forest or Gradient Boosting, that can effectively predict house prices based on the available features.
- Train the model on the prepared data, ensuring it generalizes well to unseen data. Techniques like cross-validation can be used to evaluate model performance.

➤ **Titanic Survivor Prediction:**

- Select a classification model, such as Logistic Regression or Decision Tree, that can accurately predict passenger survival based on the passenger information.
- Train the model using the Titanic dataset, again employing techniques like cross-validation to assess its generalizability.

## 5.3 MODULE 3: EXPLAINABLE AI INTEGRATION

### ➤ **XAI Integration Implementation:**

- Implement the chosen XAI techniques within our programming environment to seamlessly integrate them with the trained models. This involves using existing XAI techniques like LIME, SHAP.
- Ensure the integrated system can generate explanations alongside model predictions for both house prices and survival probabilities.

## 5.4 MODULE 4: EVALUATION AND REFINEMENT

### ➤ **Explanation Clarity and Informativeness:**

- Assess whether the explanations generated by the XAI techniques are clear, informative, and provide valuable insights into the models' reasoning.
- Evaluate the explanations for both house price predictions and survivor probabilities, ensuring they align with your expectations.

### ➤ **Model Performance:**

- Monitor the impact of XAI integration on the performance (prediction accuracy) of the underlying machine learning models.

## CHAPTER 6

### REQUIREMENTS

#### 6.1 SOFTWARE REQUIREMENTS

**Programming Environment:** Python 3 (version 3.6 or later is recommended)

➤ **Machine Learning and XAI Libraries:**

- **scikit-learn:** For building and training your machine learning models
- **LIME:** To generate localized explanations for individual house price predictions and survival probability predictions.
- **SHAP:** To provide global explanations of feature importance for both house price and survival prediction models.

➤ **Web App Development Framework:**

- **Streamlit:** To create a user-friendly web application to showcase the integrated Explainable AI models.

#### 6.2 HARDWARE REQUIREMENTS

- **Processor:** A mid-range processor with a clock speed of 2.0 GHz or higher
- **RAM:** Minimum 4GB of RAM, 8GB or more is recommended.
- **Display:** A standard resolution display of 1366x768 pixels or higher
- **Connectivity:** A stable internet connection is necessary for downloading software libraries, accessing online resources, and deploying the web application using Streamlit. A minimum speed of 10Mbps is recommended for smooth interaction.
- **Storage:** Sufficient storage space will be needed to accommodate datasets, project files, and software libraries.
- **Video Memory:** While not essential for this project, having a dedicated graphics card with at least 1GB of video memory can be beneficial for certain visualization tasks.

## CHAPTER 7

### SYSTEM DESIGN

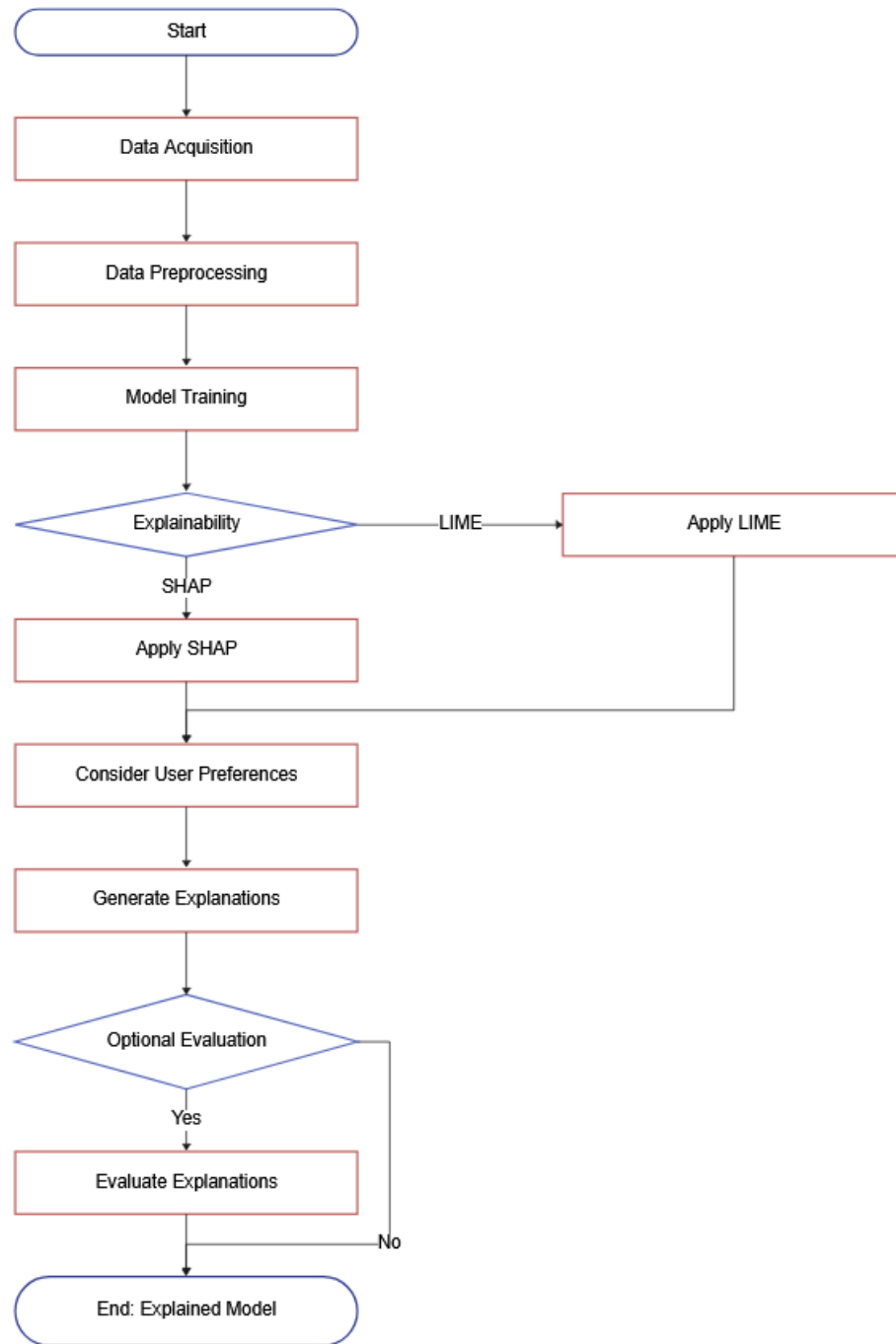


Figure 7.1 System Design Overview



## CHAPTER 8

### SAMPLE CODE

#### **main.py**

```
import streamlit as st

import streamlit.components.v1 as components

from Housing_Price_Prediction.house_price_prediction import predict_house_price

from Titanic_Survivor_Predictor.titanic_survivor_prediction import predict_survivor

import shap

import xgboost

def main():

    st.title("Explainable AI")

    st.write("Choose a prediction task:")

    choice = st.radio("", ("House Price Prediction", "Titanic Survivor Prediction"))

    if choice == "House Price Prediction":

        house_price_page()

    elif choice == "Titanic Survivor Prediction":

        titanic_survivor_page()

def house_price_page():

    st.subheader("House Price Prediction")
```

```

st.write("Enter the following details:")

lot_area = st.number_input("Lot Area")

year_built = st.number_input("Year Built")

first_floor_sf = st.number_input("1st Floor Sq Ft")

second_floor_sf = st.number_input("2nd Floor Sq Ft")

full_bath = st.number_input("Number of Full Bathrooms")

bedrooms = st.number_input("Number of Bedrooms")

total_rooms = st.number_input("Total Rooms Above Ground")

user_input = {

    "LotArea": lot_area,

    "YearBuilt": year_built,

    "1stFlrSF": first_floor_sf,

    "2ndFlrSF": second_floor_sf,

    "FullBath": full_bath,

    "BedroomAbvGr": bedrooms,

    "TotRmsAbvGrd": total_rooms

}

prediction, shap_html = predict_house_price(user_input)

st.write(f"Predicted Price: ${prediction:.2f}")

if st.button("Show SHAP Explanation"):

```

```

    st.subheader("SHAP Explanation")

    st_shap(shap_html)

def titanic_survivor_page():

    st.subheader("Titanic Survivor Prediction")

    st.write("Enter passenger details:")

    sex = st.selectbox("Sex", ["Male", "Female"])

    age = st.number_input("Age", min_value=0, max_value=150, step=1)

    fare = st.number_input("Fare", min_value=0.0, step=0.01)

    pclass = st.selectbox("Passenger Class", [1, 2, 3])

    sibsp = st.number_input("Number of Siblings/Spouses Boarded", min_value=0, step=1)

    embarked = st.selectbox("Embarked From", ["C", "Q", "S"])

    parch = st.number_input("Number of Parents/Children Boarded", min_value=0, step=1)

    user_input = {

        "Sex_female": 1 if sex == "Female" else 0,

        "Sex_male": 1 if sex == "Male" else 0,

        "Age": age,

        "Fare": fare,

        "Pclass": pclass,

        "SibSp": sibsp,

        "Embarked_C": 1 if embarked == "C" else 0,

```

```

    "Embarked_Q": 1 if embarked == "Q" else 0,

    "Embarked_S": 1 if embarked == "S" else 0,

    "Parch": parch,

}

prediction, lime_html = predict_survivor(user_input)

st.write(f"Prediction: {'Will Survive' if prediction[0] == 1 else 'Will Die'}")

if st.button("Show LIME Explanation"):

    st.subheader("LIME Explanation")

    st.components.v1.html(lime_html, height=800)

def st_shap(shap_html):

    shap_html_with_js = f"<head>{shap.getjs()}</head><body><div style='overflow: hidden;margin-top: 100px;padding: 20px;'>{shap_html}</div></body>"

    components.html(shap_html_with_js, height=700, width=1500)

if __name__ == "__main__":

    main()

```

### **house\_price\_prediction.py**

```

import pandas as pd

from sklearn.ensemble import RandomForestRegressor

import shap

def predict_house_price(user_input):

```

```

home_data = pd.read_csv('Housing_Price_Prediction/train.csv')

y = home_data['SalePrice']

features = ['LotArea', 'YearBuilt', '1stFlrSF', '2ndFlrSF', 'FullBath', 'BedroomAbvGr',
'TotRmsAbvGrd']

rf_model = RandomForestRegressor(random_state=1)

rf_model.fit(home_data[features], y)

user_data = pd.DataFrame(user_input, index=[0])

predicted_price = rf_model.predict(user_data)[0]

explainer = shap.TreeExplainer(rf_model)

shap_values = explainer.shap_values(user_data)

shap_html = shap.force_plot(explainer.expected_value, shap_values[0], user_data,
matplotlib=False)

return predicted_price, shap_html._repr_html_()

```

### **titanic\_survivor\_prediction.py**

```

import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.ensemble import RandomForestClassifier

import lime

import lime.lime_tabular

def predict_survivor(user_input):

```

```

data = pd.read_csv("Titanic_Survivor_Predictor/train.csv")

train, test = train_test_split(data, test_size=0.3, random_state=0, stratify=data['Survived'])

train = train.drop(['Name', 'Ticket', 'Cabin', 'PassengerId'], axis=1)

test = test.drop(['Name', 'Ticket', 'Cabin', 'PassengerId'], axis=1)

train_processed = pd.get_dummies(train)

test_processed = pd.get_dummies(test)

train_processed = train_processed.fillna(train_processed.mean())

test_processed = test_processed.fillna(test_processed.mean())

X_train = train_processed.drop(['Survived'], axis=1)

Y_train = train_processed['Survived']

X_test = test_processed.drop(['Survived'], axis=1)

Y_test = test_processed['Survived']

random_forest = RandomForestClassifier(n_estimators=100)

random_forest.fit(X_train, Y_train)

predict_fn_rf = lambda x: random_forest.predict_proba(x).astype(float)

X = X_train.values

explainer = lime.lime_tabular.LimeTabularExplainer(X, feature_names=X_train.columns,

                                                    class_names=['Will Die', 'Will Survive'], kernel_width=5)

preprocessed_input = preprocess_user_input(user_input)

prediction = random_forest.predict(preprocessed_input)

```

```

    explanation    =    explainer.explain_instance(preprocessed_input.values[0],    predict_fn_rf,
num_features=10)

    lime_html = explanation.as_html()

    return prediction, lime_html

def preprocess_user_input(user_input):

    user_df = pd.DataFrame([user_input])

    if 'Sex' in user_df.columns:

        user_df.drop('Sex', axis=1, inplace=True)

    user_df = user_df.reindex(columns=[ 'Pclass', 'Age','SibSp', 'Parch', 'Fare', 'Sex_female',
'Sex_male', 'Embarked_C', 'Embarked_Q', 'Embarked_S'], fill_value=0)

    return user_df

```

## CHAPTER 9

# RESULT AND OUTPUT

### Explainable AI

Choose a prediction task:

- ☒ House Price Prediction  
☐ Titanic Survivor Prediction

#### House Price Prediction

Enter the following details:

Lot Area  
0.00

Year Built  
0.00

1st Floor Sq Ft  
0.00

2nd Floor Sq Ft  
0.00

Number of Full Bathrooms  
0.00

Number of Bedrooms  
0.00

Total Rooms Above Ground  
0.00

Predicted Price: \$50984.77

Show SHAP Explanation

### Explainable AI

Choose a prediction task:

- ☐ House Price Prediction  
☒ Titanic Survivor Prediction

#### Titanic Survivor Prediction

Enter passenger details:

Sex  
Male

Age  
0

Fare  
0.00

Passenger Class  
1

Number of Siblings/Spouses Boarded  
0

Embarked from  
C

Number of Parents/Children Boarded  
0

Prediction: Will Survive

Show LIME Explanation

Figure 9.1 Home Page



## House Price Prediction

Enter the following details:

Lot Area

20000.00

- +

Year Built

2003.00

- +

1st Floor Sq Ft

215.00

- +

2nd Floor Sq Ft

210.00

- +

Number of Full Bathrooms

2.00

- +

Number of Bedrooms

5.00

- +

Total Rooms Above Ground

8.00

- +

Predicted Price: \$167599.79

Figure 9.2 House Price Prediction Model

## Titanic Survivor Prediction

Enter passenger details:

Sex

Male



Age

40



Fare

20.00



Passenger Class

1



Number of Siblings/Spouses Boarded

1



Embarked From

Q



Number of Parents/Children Boarded

2



Prediction: Will Die

Figure 9.3 Titanic Survivor Prediction Model

## House Price Prediction

Enter the following details:

Lot Area

1000000.00

Year Built

1987.00

1st Floor Sq Ft

254.00

2nd Floor Sq Ft

256.00

Number of Full Bathrooms

2.00

Number of Bedrooms

8.00

Total Rooms Above Ground

10.00

Predicted Price: \$174940.90

Show SHAP Explanation

## SHAP Explanation

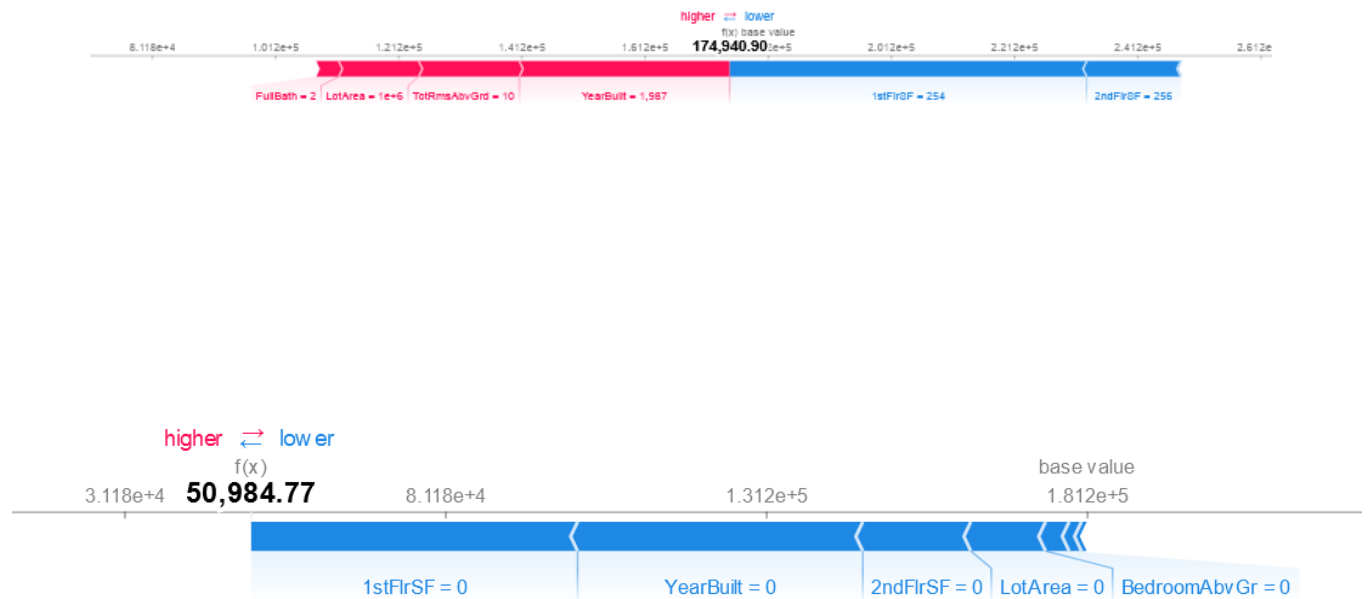


Figure 9.4 House Price Prediction With SHAP Explanation

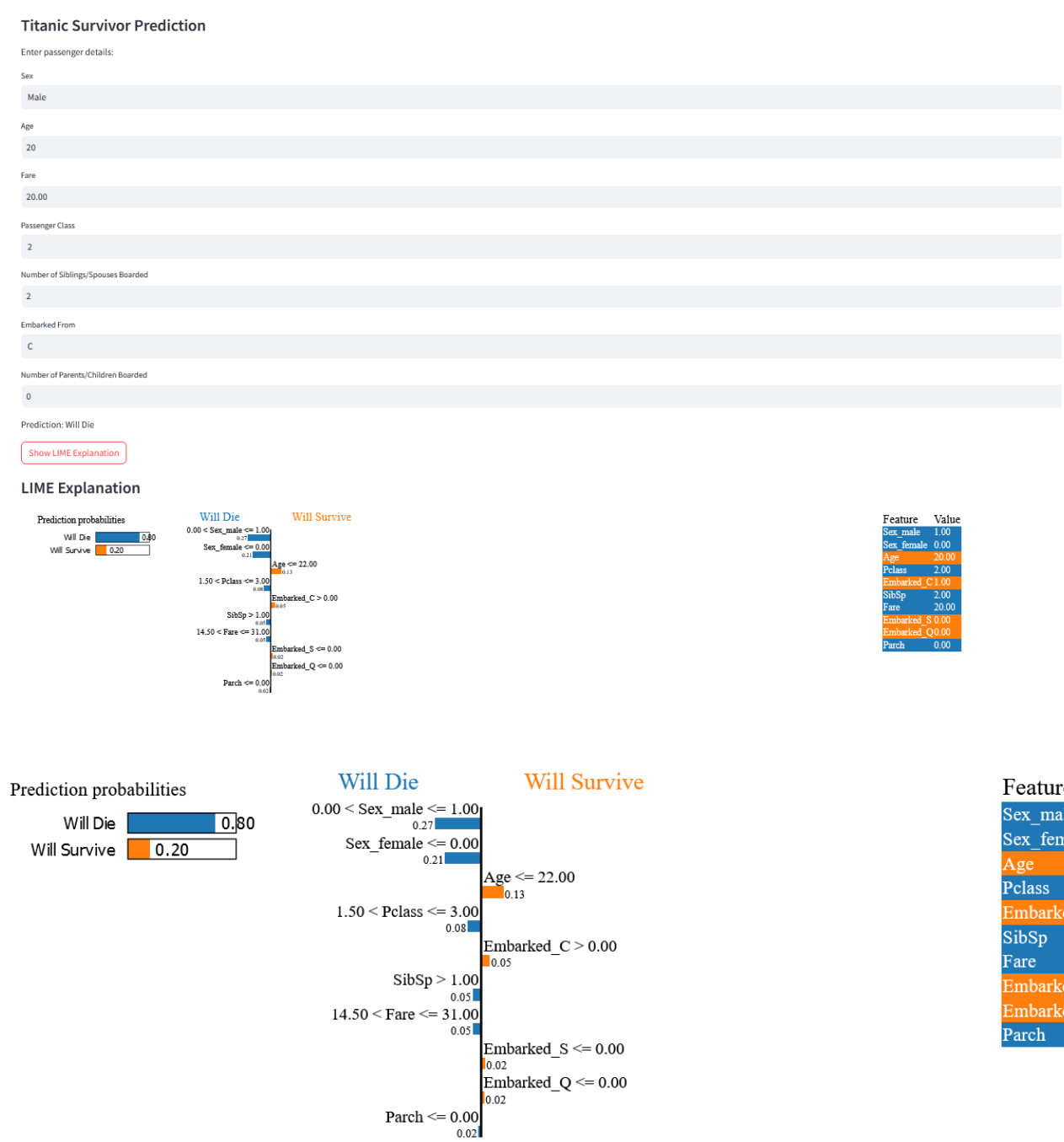


Figure 9.5 Titanic Survivor Prediction With LIME Explanation

## CHAPTER 10

### FUTURE ENHANCEMENTS

#### 10.1 EXPANDING EXPLANATIONS:

- **Advanced LIME Techniques:** Explore leveraging techniques like Anchors within LIME to not only explain individual predictions but also highlight similar data points that influenced the model's reasoning.
- **Counterfactual Explanations:** Implement techniques that generate counterfactual explanations. These explanations answer "what-if" scenarios.

#### 10.2 BEYOND REAL ESTATE AND TITANIC DATA:

The XAI integration approach outlined in this project can be applied to various machine learning domains:

- **Loan Approval Systems:** Integrate XAI techniques to explain loan approval/rejection decisions. This fosters trust and transparency for both lenders and borrowers by highlighting the factors influencing these decisions.
- **Medical Diagnosis Systems:** By integrating XAI, medical professionals can gain insights into the reasoning behind AI-powered diagnostic tools. This can enhance trust in these systems and potentially lead to improved healthcare decision-making.

#### 10.3 HUMAN-AI COLLABORATION:

As XAI techniques continue to evolve, they can facilitate a more collaborative approach between humans and AI systems:

- **Interactive Explanations:** Develop interactive interfaces that allow users to explore explanations in more detail. This can involve filtering explanations by features, visualizing data distributions, and drilling down into specific aspects of the model's reasoning.

## CHAPTER 11

### CONCLUSION

We embarked on a voyage of discovery, venturing into the enigmatic realm of black box models and illuminating the path towards a future guided by Explainable AI (XAI). We delved into the limitations of non-explainable models, exposing their opacity and susceptibility to bias. These opaque models, while capable of generating impressive results, leave us in the dark about their reasoning, hindering trust and hindering our ability to improve them.

The project journey wasn't just about the destination; it was about equipping you with the tools and knowledge to navigate the path yourself. We provided a practical roadmap, outlining the essential modules for XAI integration, the necessary hardware and software requirements, and valuable resources to guide your exploration. This empowers you to integrate XAI into your own machine learning endeavors, regardless of domain or application.

The future of XAI is brimming with possibilities. We explored the potential for future enhancements, delving into advanced XAI techniques like counterfactual explanations that allow users to explore "what-if" scenarios. We also discussed the broader application of XAI principles across various machine learning domains, from loan approval systems to medical diagnosis tools and recommender systems. As XAI becomes more sophisticated, its reach will extend to even more aspects of our lives, ensuring transparency and responsible development in a world increasingly reliant on AI.

The journey towards Explainable AI is just beginning, and this book serves as your launchpad. It equips you with the foundational knowledge and practical steps to integrate XAI into your own machine learning projects. Embrace the power of XAI, unveil the black boxes of the past, and unlock a future where AI is not just intelligent, but also responsible, transparent, and a force for good in the world. Let us embark on this exciting exploration together, for the potential of Explainable AI is truly limitless.

## CHAPTER 12

### REFERENCES

- [1] Explainable AI: DARPA's Explainable AI (XAI) Program by DARPA (Defense Advanced Research Projects Agency) [<https://www.darpa.mil/program/explainable-artificial-intelligence>](<https://www.darpa.mil/program/explainable-artificial-intelligence>)
- [2] Human-Centered Explainable AI: A Survey (2020) by Sameer Singh et al. (This paper explores the importance of user-centric explanations in XAI systems.)
- [3] scikit-learn documentation: [<https://scikit-learn.org/stable/>](<https://scikit-learn.org/stable/>) (Documentation for the scikit-learn library, commonly used for building machine learning models)
- [4] LIME (Local Interpretable Model-agnostic Explanations) documentation: [<https://github.com/marcotcr/lime>](<https://github.com/marcotcr/lime>) (Documentation for the LIME library, used for generating local explanations for individual predictions)
- [5] SHAP documentations : [<https://shap.readthedocs.io/>](<https://shap.readthedocs.io/>) (Documentation for the SHAP library, used for providing global explanations of feature importance)
- [6] Streamlit documentation: [<https://docs.streamlit.io/>](<https://docs.streamlit.io/>) (Documentation for the Streamlit library, used for creating user-friendly web applications)
- [7] Interpretable Matrix Factorization for Recommender Systems (2016) by Steffen Rendle et al. [<https://arxiv.org/pdf/2307.05680>](<https://arxiv.org/pdf/2307.05680>) (This paper explores interpretable techniques for matrix factorization, a common approach in recommender systems.)
- [8] Explainable Recommendation Systems: Survey and Challenges (2018) by Bartosz Wilk et al. [[https://arxiv.org/pdf/2202.06466.pdf?trk=public\\_post\\_comment-text](https://arxiv.org/pdf/2202.06466.pdf?trk=public_post_comment-text)]([https://arxiv.org/pdf/2202.06466.pdf?trk=public\\_post\\_comment-text](https://arxiv.org/pdf/2202.06466.pdf?trk=public_post_comment-text)) (This survey paper provides a comprehensive overview of explainability techniques in recommender systems)

[9] Towards User-Centric Explainable Recommendation Systems (2020) by Lei Xu et al. [<https://arxiv.org/abs/1804.11192>](<https://arxiv.org/abs/1804.11192>) (This paper emphasizes the importance of tailoring explanations to user needs in recommender systems)

[10] Interpretable Neural Collaborative Filtering for Explainable Recommendation Systems(2019) by Xiangnan He et al. (This paper explores interpretable deep learning techniques for recommender systems)

[11] Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow by Aurélien Géron. This book provides a comprehensive introduction to machine learning techniques using popular libraries.

[12] "Interpretable Machine Learning: A Guide for Making Black Boxes Transparent" by Christoph Molnar. This book offers a detailed exploration of interpretable machine learning approaches.

[13] "A Survey of Explainable Artificial Intelligence (XAI)" by Zachary C. Lipton. [<https://arxiv.org/pdf/1907.07374>](<https://arxiv.org/pdf/1907.07374>) (This survey paper provides a broader overview of XAI techniques and research areas)