

Project Title: Identifying and Mitigating Bias in AI Training Data

PHASE 1 - PROBLEM ANALYSIS

College Name: Maratha Mandal Engineering College

Group Members:

- Name: Mohammed Sayeel Shigganvi
CAN ID: CAN_33860603
Contribution: Abstract, problem definition, and implementation planning.
- Name: Abdussuban Shaboddin Patel
CAN ID: CAN_34000109
Contribution: Functional requirements and tools/platform selection.
- Name: Tufailahmed M Bargir CAN ID: CAN_34002247
Contribution: Non-functional requirements and model development.
- Name: Mohammed Shaibaj Shaikh CAN ID: CAN_33990553
Contribution: Visualization, reporting, and project documentation.

1. Abstract

This project focuses on developing an AI-based tool to detect and mitigate biases in training data distributions. Bias in AI training datasets can lead to skewed models that propagate unfair decisions. By leveraging AI techniques, the system will analyze datasets for biases, suggest corrective measures, and provide actionable insights through visualizations. The solution aims to enhance fairness and inclusivity in AI systems.

2. Problem Definition

Bias in AI training data leads to unfair, skewed models that perpetuate discrimination and inaccuracies in decision-making. Common issues include imbalanced datasets, representation bias, and label bias, which can affect the performance and fairness of AI systems. These biases result in unequal treatment of certain groups, reduced model accuracy, and diminished trust in AI solutions. Addressing this issue is critical to ensuring ethical, fair, and inclusive AI development and deployment.

2.1 Key Questions:

- How can biases in training data be effectively identified?
- What methodologies can mitigate these biases without compromising data integrity?
- How can the tool be integrated into existing AI workflows?

2.2 Target Users:

- Data Scientists: To ensure unbiased datasets.
- AI Developers: To integrate bias mitigation into model pipelines.
- Compliance Teams: To align AI systems with ethical guidelines.

2.3 Goal:

To develop a tool that detects and corrects biases in training datasets, ensuring fair and ethical AI systems

3. Requirements

3.1 Functional Requirements:

3.1.1 Data Ingestion and Integration:

- Support for various data formats (e.g., CSV, JSON).
- Seamless integration with data lakes and pipelines.

3.1.2 Data Preprocessing and Cleaning:

- Handle missing values and outliers.
- Normalize data for uniform representation.

3.1.3 AI-Driven Data Quality Enhancement:

- Detect statistical and algorithmic biases.

- Suggest methods to re-balance biased distributions.

3.1.4 Data Quality Monitoring:

- Continuous monitoring for new biases introduced in real-time data.

3.1.5 Report Generation and Visualization:

- Generate detailed bias analysis reports.
- Provide interactive visualizations for stakeholders.

3.2 Non-Functional Requirements:

3.2.1 Scalability:

- Handle large-scale datasets efficiently.

3.2.2 Real-time Performance:

- Analyze streaming data with minimal latency.

3.2.3 Data Security:

- Ensure compliance with data privacy standards.

3.2.4 Maintainability:

- Modular design for easy updates and integration.

3.2.5 Usability:

- User-friendly interface with minimal technical expertise required.
-

4. Tools and Platforms

4.1 Tools:

4.1.1 Data Preprocessing:

- Python, Pandas, NumPy.

4.1.2 Model Development:

- TensorFlow, PyTorch, or Scikit-learn.

4.1.3 Visualization:

- Matplotlib, Seaborn, Plotly.

4.2 Platforms:

4.2.1 Data Storage:

- AWS S3, Google Cloud Storage.

4.2.2 Data Preprocessing:

- Jupyter Notebook, Databricks.

4.2.3 Model Training & Deployment:

- AWS SageMaker, Google AI Platform.
-

5. Implementation Plan

Step 1: Data Preparation

- Gather training datasets from diverse sources.
- Preprocess data by cleaning, normalizing, and identifying preliminary biases.

Step 2: Model Development

- Develop models to detect bias using statistical and machine learning techniques.
- Implement algorithms for bias correction and re-balancing.

Step 3: Deployment

- Deploy the bias detection tool as an API for easy integration.
- Test the API with sample datasets.

Step 4: Reporting and Visualization

- Generate comprehensive reports on data biases.
- Create visual dashboards to display analysis results interactively.

6. Expected Outcomes

1. AI Model for Bias Detection and Correction:

An AI model capable of identifying and mitigating biases in training data.

2. API for Integration:

An API to integrate bias detection and correction capabilities into existing workflows.

3. Visualizations and Reports:

Interactive dashboards and detailed reports to communicate findings to stakeholders effectively.