



Data Governance and Security in AI

AI Ethics, Bias, & Transparency

Learning Objective

By the end of this lesson, learners will be able to **identify bias in AI, explain its societal impact, and evaluate methods for improving fairness and transparency in AI systems.**

Overview

Ethical AI development requires mitigating bias, ensuring transparency, and building trust in AI models. Organizations must proactively identify and address bias while implementing explainability techniques to enhance AI governance.

1. Understanding AI Bias & Fairness Challenges

Types of AI Bias

1. **Data Bias** – Occurs when training data **underrepresents certain groups** or includes historical discrimination.
2. **Model Bias** – Arises when algorithmic decisions **reinforce societal inequalities**.
3. **User Bias** – Happens when **human interactions** with AI lead to biased outcomes.

Societal Impact of AI Bias

- **Hiring Discrimination** – AI-based hiring systems may favor certain demographics over others.
- **Financial Inequality** – Loan approval models can deny loans disproportionately to minority applicants.
- **Criminal Justice Issues** – AI-driven risk assessments may overestimate recidivism rates for certain groups.

2. Transparency & Explainability in AI

What Makes AI Transparent?

- **Explainability** – Understanding how AI arrives at decisions.
- **Auditability** – The ability to inspect and validate AI models.
- **User Interpretability** – Ensuring outputs are understandable to non-technical users.

Methods for Improving AI Transparency

- **Explainable AI (XAI) Techniques:**
 - **SHAP & LIME** – Model-agnostic tools that explain AI predictions.
 - **Decision Trees & Rule-Based Models** – More interpretable alternatives to black-box models.
- **Bias Detection & Fairness Audits:**
 - Tools like **IBM AI Fairness 360** or **Google's What-If Tool**.
 - Regular fairness evaluations in model development.

Hands-On Activity: Bias & Transparency Audit

Scenario: A client in the **retail sector** is implementing an AI-powered customer segmentation model for personalized marketing. However, concerns have emerged that **certain customer groups receive fewer offers** due to algorithmic bias.

Task:

1. **Identify potential sources of bias** in the model's training data and design.
2. **Propose a fairness auditing strategy**, using XAI techniques and bias detection tools.
3. **Develop a client briefing** on ethical AI implementation strategies.

Key Takeaways

- **AI bias can emerge from data, models, or user interactions**, leading to ethical concerns.
- **Transparency & explainability** are crucial for building trust in AI.
- **XAI tools and fairness audits** help detect and mitigate bias in AI models.
- **Organizations must proactively assess and refine AI models** to ensure fairness and compliance.

[< Previous](#)

© 2025 General Assembly
[Attributions](#)