



# MODULE 5

# ETHICAL CONSIDERATIONS AND ALGORITHMIC BIAS HANDLING





# Session 2: Biasness Handling and Decision Making

Bias and Detection Techniques





- Bias in algorithms refers to systematic errors that result in unfair or prejudiced outcomes against certain groups or individuals.
- These biases can arise from the data used to train algorithms, the design choices made during development, or how algorithms are deployed.
- Detecting and mitigating bias is crucial to ensure that algorithms operate fairly and equitably.





# Types of Bias

#### 1. Data Bias:

Bias that arises from the datasets used to train algorithms. If the data is not representative of the population, the algorithm may make biased decisions.





#### Examples:

Historical Bias: If a dataset reflects past societal biases (e.g., gender discrimination in hiring), the algorithm may perpetuate these biases.

Sampling Bias: If the dataset is not diverse or representative (e.g., a facial recognition system trained only on images of lighter-skinned individuals), the algorithm may perform poorly on underrepresented groups.





# 2. Algorithmic Bias:

Bias that originates from the design and implementation of the algorithm itself. This can occur even if the data is unbiased.

# **Examples:**

Feature Selection Bias: Choosing features that inherently reflect bias (e.g., using zip codes in loan approval algorithms can indirectly introduce racial bias).





# 3. Deployment Bias:

Bias that occurs when an algorithm is used in a context different from where it was trained or intended.

# **Examples:**

Contextual Bias: Using an algorithm trained in one region or culture in another where it may not apply appropriately.

**Operational Bias:** Differences in how users interact with the system that weren't accounted for during training.





# **Detection Techniques**

# 1. Statistical Analysis:

Statistical methods can be used to identify patterns of bias in datasets or algorithmic outputs.

# Techniques:

**Descriptive Statistics:** Analyzing the distribution of data across different groups (e.g., race, gender) to identify imbalances.

**Hypothesis Testing:** Conducting tests to determine if there are significant differences in outcomes for different groups.





# 2. Fairness Metrics:

Metrics specifically designed to assess the fairness of an algorithm.

#### **Common Metrics:**

**Demographic Parity:** Ensures that the algorithm's positive outcomes are equally distributed across different demographic groups. **Equal Opportunity:** Ensures that the algorithm has equal true positive rates across groups, meaning it should be equally accurate for all





#### 3. Bias Audits:

A comprehensive review process where an algorithm is systematically tested for bias throughout its development and deployment.

# Steps:

Data Audit: Evaluating the training data for potential sources of bias.

**Model Audit:** Testing the model on different subgroups to ensure fairness.

Outcome Audit: Reviewing the impact of the algorithm's decisions in real-world settings.





# 4. Bias Mitigation Techniques:

Techniques used to reduce or eliminate bias once it has been detected.

#### Methods:

**Pre-processing:** Adjusting the data before training to remove bias.

**In-processing:** Modifying the learning algorithm to minimize bias during training.

**Post-processing:** Adjusting the model's outputs to ensure fairness after the model has been trained.



# 5. Ethical Review and Stakeholder Involvement:

Engaging stakeholders and conducting ethical reviews to identify biases that might not be apparent through technical means alone.

# Techniques:

User Testing: Involving diverse groups in testing to uncover biases that affect specific demographics.

**Ethical Committees:** Forming committees that include ethicists, social scientists, and affected communities to review algorithmic decisions.





# Fairness in Algorithmic Decision-Making



- Fairness in algorithmic decision-making refers to the principle that algorithms should make just, equitable, and unbiased decisions, ensuring that their outputs systematically disadvantage no group.
- Achieving fairness is crucial in maintaining trust, ensuring ethical use of technology, and preventing harm, particularly in critical areas like hiring, lending, law enforcement, and healthcare.





# **Understanding Fairness:**

Fairness means the decisions made by algorithms should not favor or disadvantage any individual or group based on attributes like race, gender, age, or socioeconomic status.





# Types of Fairness

# 1. Demographic Parity (Statistical Parity):

An algorithm satisfies demographic parity if the probability of a positive outcome (e.g., getting a loan or job) is the same for all demographic groups.

Example: If a hiring algorithm selects 30% of male applicants, it should also select 30% of female applicants.





# 2. Equal Opportunity:

Ensures that an algorithm has the same true positive rate (i.e., the rate at which it correctly identifies positive cases) across different groups.

Example: A healthcare algorithm should have the same accuracy in diagnosing a condition across all racial or ethnic groups.





# 3. Equalized Odds:

An algorithm meets equalized odds if it has both the same true positive rate and false positive rate across all demographic groups.

Example: A criminal justice algorithm predicting recidivism should have the same rates of false positives (wrongly predicting someone will reoffend) and true positives across all races.





# 4. Predictive Parity:

Predictive parity ensures that the predicted probability of a positive outcome (e.g., creditworthiness) is equally reliable for different groups.

Example: If a credit scoring algorithm assigns a 70% probability of repaying a loan to individuals, it should be equally accurate for men and women.





# 5. Individual Fairness:

An algorithm is individually fair if similar individuals receive similar outcomes, regardless of group membership.

Example: Two applicants with similar qualifications and experience should have the same chances of getting hired, regardless of their race or gender.





# Approaches to Ensuring Fairness

# 1. Pre-processing Techniques:

These techniques aim to remove bias from the data before it is used to train an algorithm.

#### Methods:

**Re-sampling:** Adjusting the training data to ensure equal representation of all groups.





**Re-weighting:** Assigning different weights to samples to balance representation.

**Data Augmentation:** Adding synthetic data to underrepresented groups to ensure the model learns fairly.





# 2. In-processing Techniques:

These techniques involve modifying the algorithm itself during training to reduce bias.

#### Methods:

Fairness Constraints: Incorporating fairness objectives into the algorithm's optimization process.

Adversarial Training: Training the algorithm alongside a second model designed to detect and mitigate bias.





# 3. Post-processing Techniques:

These techniques adjust the algorithm's outputs to achieve fairness after the model has been trained.

#### Methods:

**Re-ranking:** Adjusting the final rankings or scores produced by the algorithm to ensure fair outcomes. **Threshold Adjustment:** Modifying decision thresholds to ensure equal treatment across groups.





# Ethical Decision-Making Frameworks





Ethical decision-making frameworks provide structured approaches to navigate complex ethical dilemmas in coding and algorithm development.

These frameworks guide developers and organizations in making choices that align with ethical principles, ensuring that technology is used responsibly and equitably.





# Common Ethical Decision-Making Frameworks

#### 1. Utilitarianism

A consequentialist framework that focuses on the outcomes of decisions, aiming to maximize overall happiness or minimize harm.

Principle: The right action is the one that produces the greatest good for the greatest number.





#### Application:

Example: In designing a healthcare algorithm, a utilitarian approach would prioritize maximizing patient outcomes, even if it means making tradeoffs that might disadvantage some individuals.

Challenges: Utilitarianism can justify actions that harm minorities if the overall benefit is deemed greater, raising concerns about equity.





# 2. Deontology

A duty-based framework that emphasizes following ethical principles and rules, regardless of the consequences.

Principle: Certain actions are inherently right or wrong, and individuals have a duty to act accordingly.





#### Application:

Example: A deontological approach in data privacy would insist on respecting users' rights to consent and confidentiality, even if ignoring these rights might lead to beneficial outcomes.

Challenges: Deontology can lead to rigid decisions that might ignore practical outcomes, potentially resulting in harm.





#### 3. Virtue Ethics

A character-based framework that focuses on the moral character of the decision-maker rather than the rules or consequences.

Principle: Ethical decisions are those that reflect virtues such as honesty, fairness, and integrity.





#### Application:

Example: Developers might adopt a virtue ethics approach by prioritizing honesty and transparency in algorithmic development, ensuring that all stakeholders are well-informed.

Challenges: Virtue ethics can be subjective, as different cultures and individuals may have different interpretations of what constitutes a virtuous action.





# 4. Rights-Based Ethics

A framework that emphasizes the protection and respect of individual rights.

Principle: Actions are ethical if they respect the fundamental rights of individuals, such as the right to privacy, freedom, and equality.





# Application:

Example: In developing a social media platform, a rights-based approach would ensure users' rights to free expression and privacy are upheld, even when balancing these rights against other concerns like security.

Challenges: Conflicts can arise when different rights are at odds, such as privacy versus security, requiring careful balancing.





#### 5. Justice-Based Ethics

A framework that focuses on fairness, equality, and justice in decision-making.

Principle: Ethical actions are those that ensure fair treatment and equitable outcomes for all individuals.





#### Application:

Example: A justice-based approach in hiring algorithms would focus on ensuring that all candidates have an equal opportunity, regardless of their background, and that the algorithm does not perpetuate existing social inequalities.

Challenges: Ensuring fairness can be complex, especially when addressing historical inequalities or systemic bias. Care Ethics





# 6. Care Ethics

A relational framework that emphasizes the importance of relationships, empathy, and care in decision-making.

Principle: Ethical decisions are those that maintain and nurture relationships and consider the well-being of all affected parties.





#### Application:

Example: A care ethics approach in developing customer service Al would prioritize understanding and addressing the emotional and practical needs of users, ensuring the technology supports and enhances human interaction.

Challenges: Care ethics can be seen as less objective and harder to apply in situations requiring strict, rule-based decisions.





- This marks the end of Module 5.
- See you in the last Module to conclude our course by Developing practical solution to social challenges.

#### **THANK YOU**