

Evaluation Report

Analysis Date: 2025-02-22 10:03:29

Application Details

- **App Name:** Stock Research Session
- **Evaluation Mode:** batch_aggregate
- **Contract Count:** 1

Fairness Metrics

- **FTU Satisfied:** True
- **Race Words Count:** 0
- **Gender Words Count:** 0

Toxicity Metrics

- **Toxic Fraction:** 0.0
- **Max Toxicity:** 0.0
- **Toxicity Probability:** 0.0

Stereotype Metrics

- **Gender Bias Detected:** False
- **Racial Bias Detected:** False

Policy Evaluation Results

- **opa_policies\compliance\fairness\fairness.rego:** PASS
- **raw_result:** {'result': [{'expressions': [{'value': {'allow': True, 'denials': {}}, 'text': 'data.compliance.fairness.compliance_report', 'location': {'row': 1, 'col': 1}}]}]}

Evaluation Summary

Summary of Results

1. Concise Summary of the System's Fairness and Bias Metrics

- **Fairness Through Unawareness (FTU) Check:** The system passed the FTU check by not including any explicit race or gender words, indicating an attempt to avoid bias based on these attributes.
- **Toxicity Metrics:** The system shows no signs of toxicity, with a toxic fraction, maximum toxicity, and toxicity probability all at 0.0000, suggesting it handles interactions in a manner that is unlikely to offend or harm users.

2. Key Strengths in Terms of Fairness and Ethical Behavior

- **Avoidance of Explicit Bias:** The absence of race and gender words as per the FTU check indicates an effort to design an AI system that does not make decisions based on these potentially discriminatory factors.
- **Non-toxic Interaction:** The excellent toxicity metrics indicate that the system is highly unlikely to generate harmful or offensive content, contributing to a safe and respectful environment for users.

3. Areas of Concern or Potential Improvements

- **Beyond Unawareness:** While FTU is a good starting point, fairness requires more than just ignoring protected attributes. It's vital to ensure the system does not indirectly discriminate based on proxies for race, gender, or other protected characteristics.
- **Complexity of Fairness:** The FTU and toxicity metrics, while important, are relatively simplistic measures of fairness and ethical behavior. The system should also be evaluated on more nuanced dimensions of fairness (e.g., equality of opportunity, disparate impact) and tested across diverse datasets to uncover hidden biases.
- **Continuous Monitoring and Feedback:** As societal norms and definitions of fairness evolve, the system should incorporate continuous feedback and re-evaluation mechanisms to adapt to changing expectations and uncover biases that were not initially apparent.

4. Overall Assessment of the System's Suitability

- Based on the provided metrics, the system demonstrates a strong foundation in terms of avoiding explicit biases and maintaining a non-toxic environment. However, fairness and ethics in AI are multifaceted and require ongoing evaluation beyond initial metrics. The system appears suitable for deployment with the caveat that it incorporates ongoing monitoring, testing against diverse datasets, and adjustments based on feedback to ensure it remains fair and ethical in practice. Continuous improvement in response to new insights and societal changes will be key to maintaining its suitability.

Disclaimer

Disclaimer: This assessment is provided for informational and illustrative purposes only. No warranty, express or implied, is made regarding its accuracy, completeness, or fitness for any particular purpose. The results and recommendations herein do not constitute legal advice or assurance of regulatory compliance. Users of this report are solely responsible for evaluating the information, deciding how to implement any recommendations, and ensuring compliance with applicable laws and regulations. By using this report, you agree that aicertify/mantric/Principled Evolution (or any individual or organization associated with it) shall not be held liable for any direct, indirect, or consequential losses, damages, or claims arising from the use of or reliance on this

information.

CONFIDENTIAL

Disclaimer

Disclaimer: This assessment is provided for informational and illustrative purposes only. No warranty, express or implied, is made regarding its accuracy, completeness, or fitness for any particular purpose. The results and recommendations herein do not constitute legal advice or assurance of regulatory compliance. Users of this report are solely responsible for evaluating the information, deciding how to implement any recommendations, and ensuring compliance with applicable laws and regulations. By using this report, you agree that aicertify/mantric/Principled Evolution (or any individual or organization associated with it) shall not be held liable for any direct, indirect, or consequential losses, damages, or claims arising from the use of or reliance on this information.

CONFIDENTIAL