**Create a Scenario:**

A hypothetical credit card company uses AI and machine learning to detect fraud. This AI system utilizes advanced algorithms to quickly sift through massive datasets to identify irregular patterns and anomalies that may indicate fraudulent behavior. The AI system detects this fraud by comparing incoming data to a baseline of normal activity to identify deviations that could be fraudulent, which is a method known as anomaly detection. The other main ways that AI detects fraud is pattern recognition, behavioral analysis, and data verification. With pattern recognition, AI examines data to identify subtle patterns that could indicate fraud. The AI system can also analyze customer, account, and device behavior to identify patterns that could be fraudulent. This overall makes fraud detection occur faster and more accurate than traditional methods.

**Apply the AI RMF & Address Risk Management:**

AI RMF – Mapping

It is crucial to recognize the main obstacles to the implementation of AI-based fraud detection systems. There are several different challenges, and their associated risks, related to using this system. These challenges include ethical concerns related to data privacy, algorithmic biases, and system vulnerabilities. Privacy concerns are significant due to the vast amounts of personal data that the system requires for its functionality. This can pose a risk of a data breach. Additionally, AI systems can accidentally expose or misuse sensitive information. Additionally, if poor datasets are utilized for AI training this can lead to inaccurate interpretations (bias) and mistakes in decision-making for detecting fraud. Another risk of AI-based fraud detection systems could be the vulnerabilities in these systems. Threat actors could exploit these vulnerabilities to gain unauthorized access or manipulate AI algorithms to produce incorrect outcomes and affect the integrity of these systems.

AI RMF – Measure

After mapping these risks, then it is important measure the AI risks and their impacts. This involves assessing the likelihood and potential impact of identified AI risks by tracking metrics related to trustworthiness, transparency, accuracy, and validity.

AI RMF – Manage

After measuring these risks, the next step is to apply risk management processes to address them. For managing privacy concerns related to the use of this AI system, the company can implement data minimization/data anonymization to remove identifiable information from datasets used for training AI models and data encryption practices to protect against unauthorized access in case of a data breach. For algorithm bias concerns, the company can implement techniques like de-biasing data sets, using

algorithms designed to reduce bias, and implementing fairness metrics to evaluate the AI model's outputs. Lastly, the company could implement vulnerability management strategies to ensure their systems are patched to reduce the risk of threat actors gaining unauthorized access to manipulate AI algorithms to produce incorrect outcomes and affect the integrity of these systems.

**Consider Stakeholder Impact:**

I would communicate clearly with stakeholders, including customers using/signing up for credit cards about how the AI system works, its limitations, and the steps taken to address bias. I would emphasize the importance of using an AI fraud-detection system because of the advantage it carries such as its ability to detect fraud faster and more accurate than traditional methods. I would state that when issues arise such as a mistaken fraud detection, to be clear with end-users that this is something that is likely to occur with using AI but nevertheless reinforce the advantages of AI. Lastly, I will communicate that privacy was kept at the forefront while this AI system was developed prior to its deployment.

**Continuous Improvement:**

I would ensure that continuous auditing and monitoring of AI models was implemented to monitor the system's performance of detecting fraud. Specifically, I would look at the rates of false positives compared to true positives and ensure that the percentage of false positives was kept at a satisfactory standard (agreed upon by stakeholders and end-users previously). This will ensure for timely adjustments to the AI model to correct any emergent biases or unfair outcomes. This also may include retraining the model with more diverse data, altering the algorithm to correct bias, or even pausing its use if significant issues are detected.

References

International Journal of Scientific Research and Applications. (2024). *Artificial Intelligence in Fraud detection: Revolutionizing financial security* [PDF]. Retrieved from https://ijsra.net/sites/default/files/IJSRA-2024-1860.pdf

Perception Point. (n.d.). *AI security: Risks, frameworks, and best practices*. Retrieved from https://perception-point.io/guides/ai-security/ai-security-risks-frameworks-and-best-practices/

BigID. (n.d.). *5 ways generative AI empowers data security*. Retrieved from https://bigid.com/blog/5-ways-generative-ai-empowers-data-security/