

Fraud Detection

Account and Transactional Risk Scoring

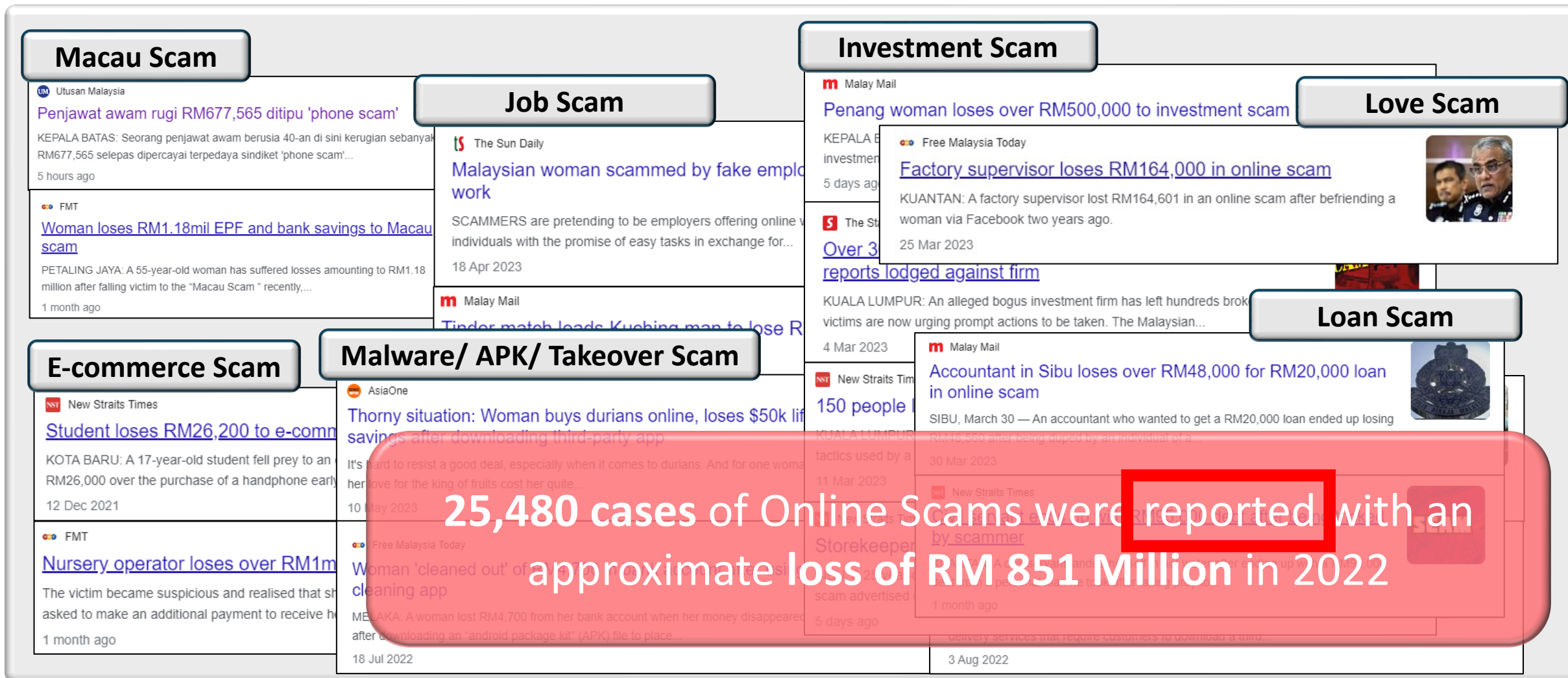
24th July 2024



What's the problem?



The ever-increasing number of fraud cases

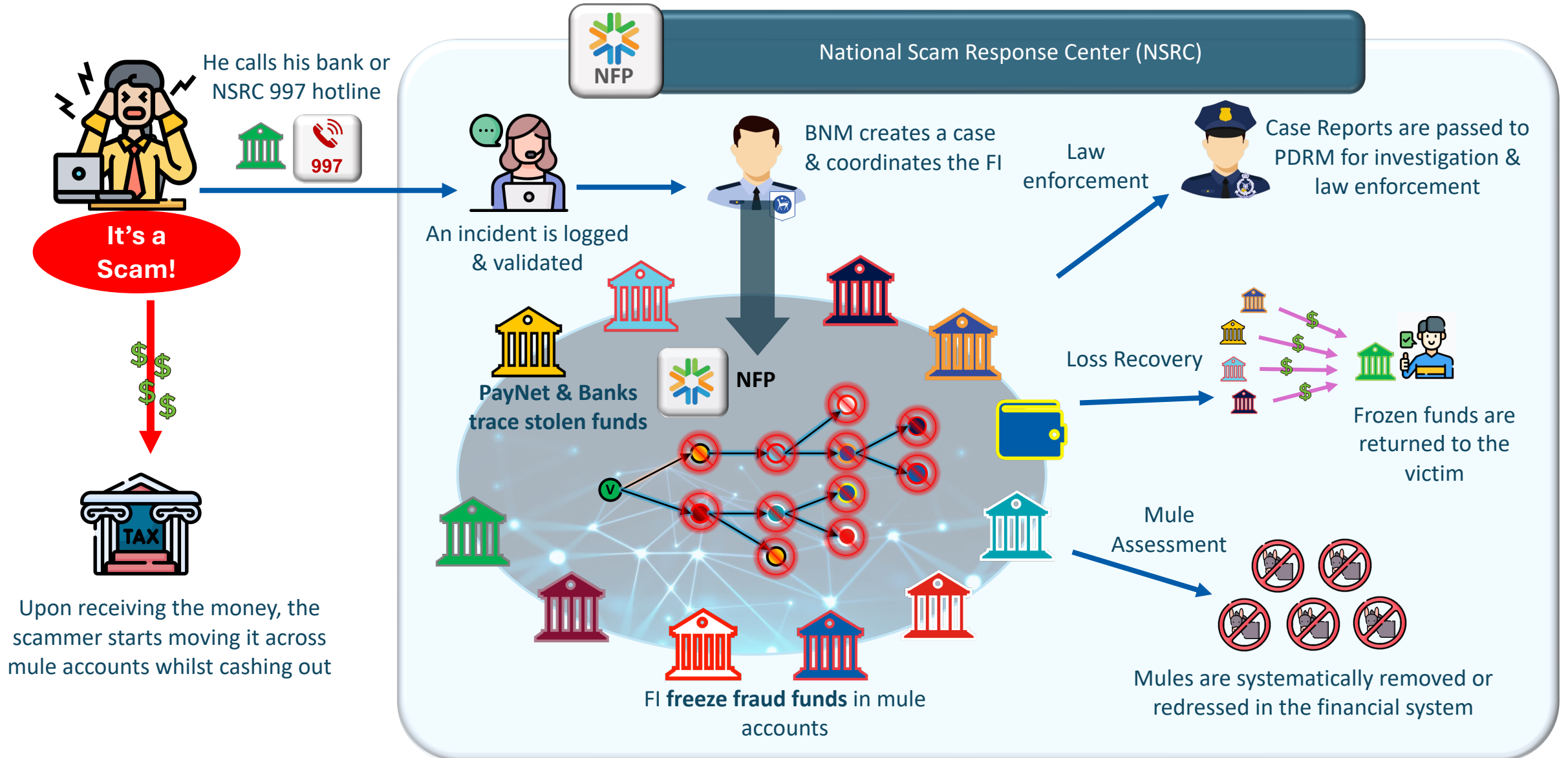




“997. What’s your emergency?”

"997. What's your emergency?"

National Scam Response Center (NSRC)





Fraud Prevention

better than cure



EMERGENCY MEETING

Fraud Prevention Challenges



Data

Data Volume

- >10 million transactions in a day
- And this is only for 2 products...

Data Cleanliness

- Missing values
- Inconsistent formats
- Reliability of information

Exploratory Data Analysis

- Basic statistics
- Distributions
- Visualize data



Labels

Extreme Imbalance

- 1000 non-fraud : 1 fraud
- Resampling / Synthetic Data / Hybrid
- Class Weights

Labels Misclassification

- K-Nearest Neighbor

Defining Fraudulent Transaction

- If an account is tagged as a mule, should all transactions by this account be considered fraudulent?
- Or, only transactions **after** the account is identified as a mule?
- **What's the difference?**



Ambiguity

Inconsistent Patterns

- Fraud patterns are quite inconsistent
- Police/Bank Officers themselves have a hard time identifying fraud transactions

Various types of fraud

- Different behaviors and patterns, for different fraud types
- How does the victim of a love scam behave prior to the scam?
- How does the feature importance differ from one type to another?

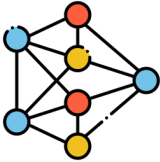
and many many more...

Fraud Prevention Model Flow



1. Feature Engineering

Arguably the most important step in machine learning!



Graph Features

- Eigenvector Centralities

Behavioral Features

- Test transaction

Interaction Features

- Features A * Features B



2. Preprocessing

Dimensionality Reduction

- e.g. Principal Component Analysis

Scaling

- Standardization/Normalization

3. Feature Selection

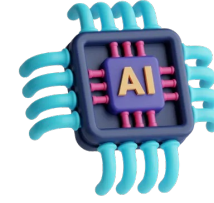
Correlation Analysis

- How does your feature affect your labels

Domain Expertise

- Selecting features based on experience

AI Bias



4. Model

XGBoost

- High Performance
- Feature Importance
- Struggle with large datasets
- Hyperparameters tuning



- Complex Patterns
- Feature Learning
- Black Box
- Very data hungry!



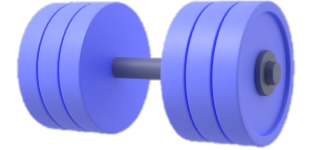
Random Forest

- Relatively little tuning required
- Imbalanced data
- Heavy (Size and Compute)

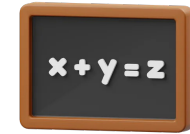


Isolation Forest

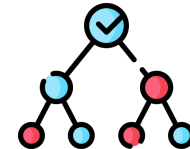
- Unsupervised Learning
- Good explainability!
- Limited to Anomaly Detection



5. Training



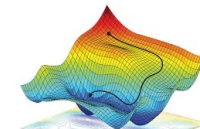
The slight changes to the weights and biases of a formula



Addition of trees based on optimal split point



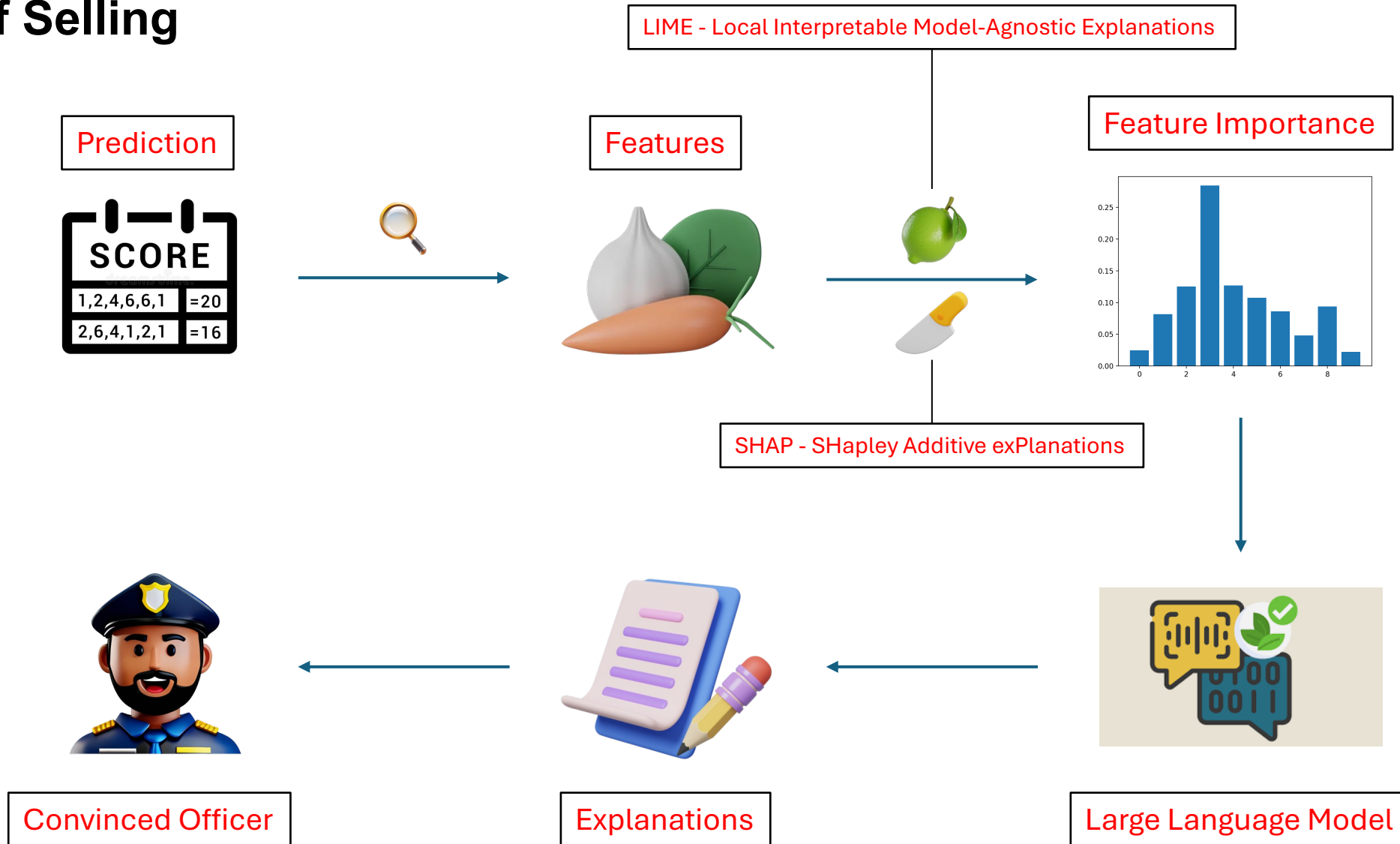
Cost Function (except Iso. F.)



Optimization Function (except Iso. F.)



Explainability





Thanks for listening!