# Mathematical Framework of Academic Misconduct Detection: Comprehensive Analysis

Leon Nduati

May 2025

## 1 Feature Space and Data Representation

### Mathematical Definition

Let us define our student data as a collection of $n$ samples $\{(x_i, y_i)\}_{i=1}^n$ where:

- $x_i \in \mathbb{R}^d$ represents the feature vector for student $i$, with $d = 9$ features.

- $y_i \in \{0, 1\}$ represents the class label where 1 indicates academic misconduct and 0 indicates normal behavior.

The feature space $\mathcal{X} \subset \mathbb{R}^d$ encompasses all possible combinations of student performance metrics. Each sample $x_i$ is represented as:

$$x_i = \begin{bmatrix} z_{c,i} \\ z_{e,i} \\ z_{d,i} \\ \sigma_i \\ t_i \\ p_i \\ v_i \\ h_i \\ a_i \end{bmatrix}$$

### Explanation

This vector representation captures the multi-dimensional nature of student performance. Each dimension represents a specific aspect of academic behavior:

- The z-scores $z_{c,i}$ and $z_{e,i}$ standardize raw performance metrics, enabling fair comparison across different assessments.

- The difference between coursework and exam performance $z_{d,i}$ highlights potential inconsistencies.

- The remaining features capture temporal patterns $t_i$, peer-relative performance $p_i$, consistency across subjects $v_i$, historical trends $h_i$, and overall statistical anomalies $a_i$.

This comprehensive feature set was carefully designed to capture various cheating behaviors, including coursework advantage (plagiarism), exam advantage (having answers in advance), and suspicious consistency patterns.

# 2 Feature Engineering Transformations

## 2.1 Z-Score Transformation

**Mathematical Definition**

For a given raw score $s_{j,i}$ where $j \in \{c, e\}$ indicating coursework or exam:

$$z_{j,i} = \frac{s_{j,i} - \mu_j}{\sigma_j}$$

where $\mu_j = \frac{1}{n} \sum_{i=1}^{n} s_{j,i}$ and $\sigma_j = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (s_{j,i} - \mu_j)^2}$.

**Explanation**

Z-score normalization serves multiple purposes in this context:

1. **Standardization**: Converts raw scores to standard deviations from the mean, allowing fair comparison.

2. **Distribution normalization**: Creates approximately normal distributions under typical conditions.

3. **Outlier highlighting**: Makes unusual performances ($\pm 2$ or more standard deviations) immediately apparent.

The z-score difference $z_{d,i}$ is particularly important as it quantifies the inconsistency between coursework and exam performance. In legitimate academic scenarios, this difference typically follows a normal distribution centered near zero. Significant deviations suggest potential misconduct.

## 2.2 Subject Variation Metric

**Mathematical Definition**

Given scores across $m$ subjects $\{s_{i,1}, s_{i,2}, \ldots, s_{i,m}\}$ for student $i$:

$$v_i = \sqrt{\frac{1}{m} \sum_{j=1}^{m} (s_{i,j} - \overline{s}_i)^2}$$

where $\overline{s}_i = \frac{1}{m} \sum_{j=1}^{m} s_{i,j}$ is the mean score across subjects.

**Explanation**

The subject variation metric $v_i$ is essentially the standard deviation of a student's performance across different subjects. This metric captures a key behavioral pattern:

- **Normal academic behavior**: Students typically show natural variation across subjects based on aptitude and interest.

- **Cheating behavior**: Students engaging in misconduct often show unnaturally consistent performance, particularly in subjects where they cheat.

- **Quantification**: Lower values of $v_i$ than expected (given a student's overall performance) can signal suspicious consistency.

This metric is especially effective at identifying students who might be obtaining unauthorized assistance across multiple subjects, as their performance becomes artificially uniform.

## 2.3   Isolation Forest Anomaly Score

**Mathematical Definition**

The anomaly score $a_i$ for student $i$ is defined as:

$$a_i = 2^{-\frac{E[h(x_i)]}{c(n)}}$$

where:

- $h(x_i)$ is the path length for sample $x_i$ (number of edges traversed in isolation tree).

- $E[h(x_i)]$ is the average path length across the forest.

- $c(n) = 2H(n-1) - \frac{2(n-1)}{n}$ where $H(i) \approx \ln(i) + 0.5772156649$ (Euler's constant).

**Explanation**

The Isolation Forest algorithm is particularly suited for detecting outliers in academic data because:

1. **Unsupervised detection**: It identifies outliers without requiring labeled cheating examples.

2. **Efficiency**: It isolates anomalies quickly by partitioning the feature space.

3. **Path length principle**: Anomalous points require fewer partitions to isolate, resulting in shorter paths.

The normalization by $c(n)$ (the average path length in a binary search tree) calibrates the score to a $[0, 1]$ range, where values closer to 1 indicate greater anomaly likelihood. This score serves as a robust, model-agnostic indicator of unusual behavior patterns that complements the other engineered features.

In the context of our model, this feature acts as an ensemble within an ensemble, providing an orthogonal view of the data that might catch anomalies missed by other features.

# 3 Model Training and Optimization

## 3.1 SMOTE Balancing

### 3.1.1 Mathematical Definition

The SMOTE algorithm creates synthetic samples $\tilde{x}$ for the minority class:

$$\tilde{x} = x_i + \lambda \cdot (x_j - x_i)$$

where:

- $x_i$ is a sample from the minority class

- $x_j$ is one of the $k$-nearest neighbors of $x_i$ in the minority class

- $\lambda \in [0, 1]$ is a random number

### 3.1.2 Explanation

SMOTE (Synthetic Minority Oversampling Technique) addresses class imbalance by generating synthetic examples of the minority class through interpolation. This approach:

1. **Avoids duplicating existing samples**: Unlike simple oversampling, it creates new, unique examples.

2. **Preserves feature relationships**: Maintains the underlying patterns of cheating behavior.

3. **Expands the minority class representation**: Allows the model to better learn decision boundaries.

By generating synthetic samples along line segments connecting real minority class samples, SMOTE creates realistic examples of potential cheating behaviors that help the model generalize better.

## 3.2 StandardScaler Transformation

### 3.2.1 Mathematical Definition

For each feature $f$, the standardized value is:

$$\tilde{x}_{i,f} = \frac{x_{i,f} - \mu_f}{\sigma_f}$$

where $\mu_f$ and $\sigma_f$ are the mean and standard deviation of feature $f$ across all samples.

### 3.2.2 Explanation

Feature standardization ensures that all features contribute equitably to the model regardless of their original scales. This preprocessing step:

1. **Equalizes feature importance**: Prevents features with larger numeric ranges from dominating.

2. **Improves convergence**: Facilitates optimization during model training.

3. **Maintains interpretability**: Allows comparison of feature importances on a common scale.

In the academic misconduct detection context, this ensures that metrics like historical trends and peer comparisons are weighted appropriately relative to z-scores and time-based features.

## 3.3 Grid Search Optimization

### 3.3.1 Mathematical Definition

The optimal hyperparameters $\theta^*$ are determined by:

$$\theta^* = \arg\max_{\theta \in \Theta} \mathrm{F1}_{\mathrm{CV}}(\theta)$$

where:

- $\Theta$ is the hyperparameter search space

- $\mathrm{F1}_{\mathrm{CV}}(\theta)$ is the average F1-score across cross-validation folds

### 3.3.2 Explanation

Grid search explores a predefined hyperparameter space to find the configuration that maximizes model performance. For academic misconduct detection:

1. **F1-score optimization**: Balances precision and recall, crucial for fair assessment.

2. **Cross-validation robustness**: Ensures performance generalizes across different data subsets.

3. **Domain-specific tuning**: Focuses on parameters most relevant to capturing cheating patterns.

The optimal configuration found (max_depth=10, max_features='sqrt', min_samples_leaf=1, min_samples_split=2, n_estimators=100) balances model complexity with generalization ability, avoiding both underfitting and overfitting.

## 3.4 Threshold Optimization

### 3.4.1 Mathematical Definition

The optimal classification threshold $\tau^*$ is determined by:

$$\tau^* = \arg\max_{\tau \in [0,1]} F_\beta(\tau)$$

where $F_\beta$ is the F-beta score:

$$F_\beta = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{\beta^2 \cdot \text{precision} + \text{recall}}$$

For $\beta = 1$ (F1-score), precision and recall are weighted equally. The optimal threshold found is $\tau^* \approx 0.58$.

### 3.4.2 Explanation

Threshold optimization is particularly important in academic misconduct detection due to the asymmetric consequences of errors:

1. **Precision emphasis**: False accusations can severely impact students' academic careers.

2. **Recall importance**: Missing actual misconduct undermines academic integrity.

3. **Institutional policy alignment**: The threshold can be adjusted to reflect specific institutional priorities.

The threshold of 0.58 (rather than the default 0.5) indicates a slight preference for precision over recall in the model, requiring stronger evidence before flagging a student. This aligns with the ethical principle of "innocent until proven guilty" while still maintaining effective detection capabilities.

## 3.5 Feature Importance Quantification

### 3.5.1 Mathematical Definition

The importance $I_f$ of feature $f$ is calculated as:

$$I_f = \frac{1}{B} \sum_{b=1}^{B} \sum_{n \in T_b} \Delta I(n, f) \cdot \frac{|S_n|}{|S_{\text{root}}|}$$

where:

- $n$ iterates over all nodes in tree $T_b$ that split on feature $f$

- $\Delta I(n, f)$ is the impurity decrease at node $n$

- $|S_n|$ is the number of samples at node $n$

- $|S_{\text{root}}|$ is the number of samples at the root node

### 3.5.2 Explanation

Feature importance provides crucial insights into which behavioral patterns most strongly indicate academic misconduct. The results from the model show:

1. **Historical trend dominance (32.7%)**: Sudden improvements in performance are the strongest indicators.

2. **Score variance significance (26.4%)**: Unnatural consistency across assessments is highly suspicious.

3. **Subject variation relevance (12.8%)**: Unusual patterns across different subjects provide important signals.

4. **Temporal cues (9.9%)**: Unusual exam timing patterns contribute meaningful information.

5. **Isolation Forest integration (7.8%)**: The anomaly score adds complementary perspective.

These importance values align with educational research on academic integrity, where dramatic performance changes and suspicious consistency are well-documented indicators of potential misconduct.

## 3.6 Evaluation Metrics

### 3.6.1 Confusion Matrix

### 3.6.2 Mathematical Definition

Define:

- True Positives (TP): $\sum_{i=1}^{n} \mathbb{I}(y_i = 1 \text{ and } \hat{y}_i = 1)$

- False Positives (FP): $\sum_{i=1}^{n} \mathbb{I}(y_i = 0 \text{ and } \hat{y}_i = 1)$

- False Negatives (FN): $\sum_{i=1}^{n} \mathbb{I}(y_i = 1 \text{ and } \hat{y}_i = 0)$

- True Negatives (TN): $\sum_{i=1}^{n} \mathbb{I}(y_i = 0 \text{ and } \hat{y}_i = 0)$

### 3.6.3 Explanation

The confusion matrix provides a complete picture of model performance across all possible prediction outcomes:

- TP (27): Correctly identified instances of academic misconduct

- FP (2): Honest students incorrectly flagged as potential cheaters

- FN (4): Missed instances of academic misconduct

- TN (167): Correctly identified honest academic behavior

The very low false positive rate (2 out of 169 honest students, or approximately 1.2%) is particularly important given the serious implications of misconduct allegations.

### 3.6.4   Performance Metrics

### 3.6.5   Mathematical Definition

- Precision: $P = \frac{TP}{TP+FP} = \frac{27}{27+2} = 0.93$

- Recall: $R = \frac{TP}{TP+FN} = \frac{27}{27+4} = 0.87$

- F1-Score: $F_1 = \frac{2 \cdot P \cdot R}{P+R} = \frac{2 \cdot 0.93 \cdot 0.87}{0.93+0.87} = 0.90$

- Accuracy: $ACC = \frac{TP+TN}{TP+TN+FP+FN} = \frac{27+167}{27+167+2+4} = 0.97$

### 3.6.6   Explanation

These metrics quantify different aspects of model performance:

1. **Precision (93%)**: High precision means that when the model flags a student, it's rarely a false alarm.

2. **Recall (87%)**: Strong recall indicates that most misconduct cases are successfully identified.

3. **F1-Score (90%)**: The harmonic mean shows that the model balances precision and recall effectively.

4. **Accuracy (97%)**: Overall, the model correctly classifies the vast majority of cases.

These metrics demonstrate that the model achieves its dual objectives: identifying academic misconduct while minimizing false accusations.

# 4   Probabilistic Interpretation

## 4.1   Mathematical Definition

The final model can be interpreted as estimating the posterior probability:

$$P(y = 1|x) = P(x \in \mathcal{C}_1|x)$$

where $\mathcal{C}_1$ represents the set of feature vectors associated with academic misconduct.
For decision-making, we use the threshold rule:

$$\hat{y} = \begin{cases} 1 & \text{if } P(y = 1|x) > \tau^* \\ 0 & \text{otherwise} \end{cases}$$

## 4.2 Explanation

This probabilistic framework provides several advantages for academic integrity applications:

1. **Uncertainty quantification:** Probabilities express the strength of evidence rather than binary judgments.

2. **Risk control:** Institutions can adjust thresholds based on their tolerance for different error types.

3. **Transparency:** Probabilistic outputs can be explained to stakeholders as likelihood estimates.

By presenting results as probabilities rather than definitive labels, the system acknowledges inherent uncertainty and encourages appropriate human judgment in the review process.

# 5 Bayesian Risk Minimization

## 5.1 Mathematical Definition

The threshold $\tau^*$ can also be viewed as minimizing the Bayes risk:

$$R(h) = \mathbb{E}_{x,y}[L(y, h(x))]$$

where $L$ is a cost-sensitive loss function:

$$L(y, \hat{y}) = \begin{cases} 0 & \text{if } y = \hat{y} \\ C_{FP} & \text{if } y = 0, \hat{y} = 1 \\ C_{FN} & \text{if } y = 1, \hat{y} = 0 \end{cases}$$

The optimal threshold in this framework becomes:

$$\tau^* = \frac{C_{FP}}{C_{FP} + C_{FN}}$$

where $C_{FP}$ and $C_{FN}$ are the costs associated with false positives and false negatives, respectively.

## 5.2 Explanation

The Bayesian risk minimization perspective formalizes the ethical trade-offs inherent in academic misconduct detection:

1. **Cost asymmetry:** False accusations typically carry higher costs than missed detections.

2. **Institutional values:** The relative costs reflect institutional priorities regarding academic integrity.

3. **Formal optimization:** Provides mathematical justification for threshold selection.

With the optimal threshold of approximately 0.58, we can infer that the implicit cost ratio in our model is $\frac{C_{FP}}{C_{FN}} \approx \frac{0.58}{0.42} \approx 1.38$, suggesting that false positives are considered about 38% more costly than false negatives.

# 6 Feature Space Visualization

## 6.1 Mathematical Definition

The high-dimensional feature space can be visualized using t-SNE projection:

$$z_i = t\text{-SNE}(x_i)$$

where $z_i \in \mathbb{R}^2$ is a two-dimensional representation of $x_i$ preserving the local neighborhood structure, enabling visualization of the separation between classes.

## 6.2 Explanation

Visualizing the feature space provides several insights:

1. **Cluster identification:** Reveals natural groupings of similar academic behaviors.

2. **Boundary visualization:** Shows how the model separates legitimate and suspicious behaviors.

3. **Outlier detection:** Highlights unusual cases for further investigation.

T-SNE is particularly useful here as it preserves local structure, allowing us to see how similar cheating strategies cluster together in the reduced dimensionality space.

# 7 Additional Considerations (Extension)

## 7.1 Permutation Feature Importance

### 7.1.1 Mathematical Definition

For a feature $f$ and performance metric $M$, the permutation importance is:

$$I_f^{perm} = M(y, \hat{y}) - M(y, \hat{y}_f^{perm})$$

where $\hat{y}_f^{perm}$ represents predictions after randomly permuting feature $f$.

### 7.1.2 Explanation

Permutation importance offers a model-agnostic alternative to the built-in Random Forest importance metric:

1. **Causal perspective:** Measures the performance drop when feature information is destroyed.

2. **Correlation robustness:** Less influenced by correlations between features.

3. **Direct interpretation:** Directly tied to model performance rather than tree structure.

This approach could provide additional validation of the feature importance results, particularly for features that might be correlated like z-score differences and historical trends.

## 7.2 Calibration Analysis

### 7.2.1 Mathematical Definition

For a well-calibrated model, the following should hold for all probability values $p$:

$$P(y = 1|P(y = 1|x) = p) \approx p$$

This can be evaluated using calibration curves plotting predicted probabilities against observed frequencies.

### 7.2.2 Explanation

Probability calibration is crucial for decision-making in high-stakes contexts like academic integrity:

1. **Reliability assessment:** Ensures probabilistic outputs reflect true likelihoods.

2. **Confidence interpretation:** Allows stakeholders to interpret probability scores correctly.

3. **Fair comparison:** Enables consistent thresholds across different contexts.

Techniques like Platt scaling or isotonic regression could be applied if calibration analysis reveals systematic over- or under-confidence in probability estimates.

## 7.3 Local Explainability with SHAP Values

### 7.3.1 Mathematical Definition

The SHAP value for feature $f$ and instance $x_i$ is defined as:

$$\phi_f(x_i) = \sum_{S \subseteq F \setminus \{f\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} [f_x(S \cup \{f\}) - f_x(S)]$$

where $F$ is the set of all features and $f_x(S)$ is the prediction for instance $x_i$ using only features in subset $S$.

### 7.3.2 Explanation

While global feature importance provides overall insights, SHAP values offer instance-specific explanations:

1. **Case-specific justification:** Explains precisely why a particular student was flagged.

2. **Transparent decision-making:** Enables clear communication with stakeholders.

3. **Pattern identification:** Helps identify new or evolving cheating strategies.

This local explainability would be particularly valuable for human reviewers making final decisions about flagged cases.

# 8 Conclusion

This comprehensive mathematical framework establishes a rigorous foundation for academic misconduct detection using machine learning. By formalizing each component—from feature engineering to model training and evaluation—we ensure that the system is both technically sound and ethically responsible.

The demonstrated performance metrics (93% precision, 87% recall, 90% F1-score) validate that this approach effectively balances the competing objectives of identifying academic misconduct while minimizing false accusations. The mathematical formulation further provides transparency into the decision-making process, which is essential for fair and responsible application in educational settings.