# Literature Review

**Introduction**

With the rise of spam and phishing emails, there is a significant threat to global security and the economy. Spam emails typically contain marketing content, while phishing emails aim to deceive recipients into divulging personal information, leading to identity theft. Detecting these malicious emails is crucial, and machine learning has become a powerful tool in this domain. This literature review focuses on the applicability of Support Vector Machine (SVM), Random Forest (RF), Naive Bayes (NB), and Artificial Neural Networks (ANN) in phishing email detection, summarizing their strengths and challenges.

**Support Vector Machine (SVM)**

Support Vector Machine is robust for high-dimensional data, suitable for text classification. SVM constructs a hyperplane in a multidimensional space to separate different classes. It can handle linear and non-linear data by using kernel functions. SVMs have shown high accuracy in spam and phishing detection tasks but can be computationally intensive, especially with large datasets.

- **Advantages**:
    - Handles high-dimensional data well.
    - Can address non-linear classification problems.
    - Generally achieves high accuracy.
- **Disadvantages**:
    - Long training time for large datasets.
    - Requires tuning of parameters like kernel functions and regularization parameters.

**Random Forest (RF)**

Random Forest is an ensemble learning method that constructs multiple decision trees during training and outputs the mode of the classes for classification tasks. It is highly effective for spam detection due to its robustness and ability to handle overfitting. RF evaluates various features such as word frequency, sender address, and the presence of URLs, making it versatile and accurate. However, its complexity and computational requirements can be a drawback for very large datasets.

- **Advantages**:
    - Robust and accurate.
    - Can handle overfitting well.
    - Provides good feature importance explanations.

- **Disadvantages**:
  - Computationally intensive on large datasets.
  - May struggle with very high-dimensional sparse data.

**Naive Bayes (NB)**

Naive Bayes is a probabilistic classifier based on Bayes' theorem, assuming feature independence. It is simple, efficient, and performs well on small datasets, making it a popular choice for text classification tasks like spam detection. The algorithm calculates the probability of an email being spam or phishing based on word frequencies and other textual features.

- **Advantages**:
  - Simple and efficient.
  - Fast training and prediction times.
  - Performs well with small datasets.
- **Disadvantages**:
  - Assumes feature independence, which may not hold in practice.
  - May struggle with complex patterns where feature dependencies are significant.

**Artificial Neural Networks (ANN)**

Artificial Neural Networks, including architectures like Long Short-Term Memory (LSTM), have gained popularity for their ability to model complex, non-linear relationships in data. ANN and LSTM can capture patterns in email content and structure, making them suitable for phishing detection. LSTMs, in particular, are effective for sequential data and can identify temporal dependencies in emails.

- **Advantages**:
  - Strong ability for complex pattern recognition.
  - Suitable for sequential and text data.
  - Can achieve high accuracy.
- **Disadvantages**:
  - Requires substantial computational resources.
  - High model complexity with lower interpretability.
  - Needs large datasets for training.

## Comprehensive Consideration

For the Capstone project on phishing email detection, a balanced approach combining simplicity, efficiency, and robustness is essential. Based on the strengths and weaknesses of the reviewed algorithms, the following steps are recommended:

1. **Initial Exploration and Baseline Model**:

Use Naive Bayes (NB) as a baseline model to quickly establish a preliminary phishing email detection system. This approach allows for fast implementation and provides a benchmark for comparison.

2. **Feature Engineering and Advanced Models**:

Utilize TF-IDF or word embeddings (like Word2Vec or GloVe) for feature extraction to capture the semantic meaning of words in emails.

Employ Random Forest (RF) for its robustness and ability to handle diverse features, providing a more accurate and stable model.

Implement Support Vector Machine (SVM) to leverage its high accuracy in handling high-dimensional data, suitable for the textual nature of emails.

3. **Model Optimization and Ensemble**:

Combine the strengths of different algorithms by using ensemble methods. For instance, a voting ensemble of NB, SVM, and RF can be used to improve overall performance by leveraging the strengths of each model.

Perform hyperparameter tuning and cross-validation to optimize model performance and ensure generalizability.

4. **Evaluation and Selection**:

Evaluate model performance using multiple metrics, including accuracy, precision, recall, and F1 score. This comprehensive evaluation ensures that the chosen model balances the detection of both phishing and legitimate emails.

Select the most suitable model based on the specific requirements of the application, considering factors such as computational efficiency and the need for real-time detection.

## Conclusion

In summary, the Capstone project can leverage the strengths of Naive Bayes, Support Vector Machine, Random Forest, and Artificial Neural Networks to build an effective phishing email detection system. By integrating these algorithms and employing robust feature engineering and model optimization techniques, the project aims to achieve high accuracy and reliability in detecting phishing emails, enhancing email security and user protection.