# Final Project Mining Project Report

**Name:** Tejas Belakavadi Kemparaju
**NJIT UCID:** tb389
**Email Address:** tb389@njit.edu

**Date:** 10/10/2024
**Professor:** Yasser Abduallah
**Course:** CS 634101 Data Mining
**Github Repo**: Link

## Abstract

This project investigates supervised data mining techniques with a focus on binary classification. Three classification algorithms—Random Forest, Support Vector Machine (SVM), and Long Short-Term Memory (LSTM)—are implemented using Python libraries and applied to the Breast Cancer Wisconsin (Diagnostic) dataset. The project evaluates classification performance through metrics such as Accuracy, Precision, Recall, False Positive Rate (FPR), and False Negative Rate (FNR). Utilizing a 10-fold cross-validation approach, this study compares the efficiency and effectiveness of the algorithms, offering insights into their strengths and limitations. Results are presented in tabular form for detailed visualization.

## Introduction

Binary classification is a fundamental task in supervised learning, particularly useful for decision-making in medical, financial, and industrial domains. This project explores three distinct algorithms for binary classification:

1. **Random Forest**: A powerful ensemble learning method based on decision trees.
2. **Support Vector Machine (SVM)**: A kernel-based approach that aims to find the optimal hyperplane for classification.
3. **Long Short-Term Memory (LSTM)**: A deep learning architecture suited for sequential data but adapted here for binary classification.

The Breast Cancer Wisconsin (Diagnostic) dataset is chosen for this study due to its real-world applicability in predicting malignant or benign tumors. By leveraging existing Python libraries for algorithm implementation and manually computing performance metrics, this project offers a hands-on understanding of classification techniques and their evaluation.

## Metrics

The performance of each algorithm is evaluated using the following metrics:

- **True Positive (TP)**: Correctly predicted positive cases.
- **True Negative (TN)**: Correctly predicted negative cases.
- **False Positive (FP)**: Incorrectly predicted positive cases.
- **False Negative (FN)**: Incorrectly predicted negative cases.
- **Accuracy**: (TP + TN) / (TP + TN + FP + FN).
- **Precision**: TP / (TP + FP).
- **Recall**: TP / (TP + FN).
- **False Positive Rate (FPR)**: FP / (FP + TN).
- **False Negative Rate (FNR)**: FN / (FN + TP).

# Project Workflow

## Data Preparation and Loading

**The Dataset is in .data format and the entire code is built on it. Please use the same data set.**

1. **Dataset**:
   - The Breast Cancer Wisconsin (Diagnostic) dataset was downloaded from the UCI Machine Learning Repository.
   - Features include 30 numeric attributes representing tumor characteristics.
   - The target variable indicates whether a tumor is malignant (M) or benign (B).
2. **Preprocessing**:
   - Removed irrelevant columns (e.g., ID).
   - Encoded the target variable: Benign (B) = 0, Malignant (M) = 1.
   - Split the dataset into training (80%) and testing (20%) sets.
3. **Cross-Validation**:
   - Implemented 10-fold cross-validation to ensure robust evaluation.

## User Interaction

A user-friendly interface was developed to set support and confidence thresholds. Input validation ensures proper ranges before proceeding with analysis.

# Algorithm Implementation

## Random Forest

- Implemented using `sklearn.ensemble.RandomForestClassifier`.
- Built a forest of decision trees with randomized node splitting.
- Evaluated predictions for each fold using manually computed metrics.

## Support Vector Machine (SVM)

- Used `sklearn.svm.SVC` with a linear kernel.
- Aimed to find the hyperplane that maximizes the margin between classes.
- Applied 10-fold cross-validation for performance evaluation.

## Long Short-Term Memory (LSTM)

- Adapted for binary classification using the `tensorflow.keras` library.
- Built a sequential model with one LSTM layer and a dense output layer.
- Binary cross-entropy was used as the loss function, and Adam optimizer was used.
- Reshaped input data to match LSTM requirements and trained the model for 3 epochs per fold.

# Performance Metrics

## Random Forest

Performance metrics (average across 10 folds):

| Metric | Value |
|---|---|
| Accuracy | 96.5% |
| Precision | 94.8% |
| Recall | 95.3% |
| FPR | 3.2% |
| FNR | 4.7% |

## SVM

Performance metrics (average across 10 folds):

| Metric | Value |
|---|---|
| Accuracy | 95.7% |
| Precision | 94.0% |
| Recall | 94.5% |
| FPR | 4.0% |
| FNR | 5.5% |

## LSTM

Performance metrics (average across 10 folds):

| Metric | Value |
|---|---|
| Accuracy | 93.2% |
| Precision | 91.0% |
| Recall | 92.1% |
| FPR | 6.8% |
| FNR | 7.9% |

**Summary Comparison**

| Algorithm | Accuracy | Precision | Recall | FPR | FNR |
|---|---|---|---|---|---|
| Random Forest | 96.5% | 94.8% | 95.3% | 3.2% | 4.7% |
| SVM | 95.7% | 94.0% | 94.5% | 4.0% | 5.5% |
| LSTM | 93.2% | 91.0% | 92.1% | 6.8% | 7.9% |

# Discussion

The Random Forest algorithm demonstrated the best overall performance, achieving the highest accuracy and the lowest FPR and FNR. This is attributed to its ensemble nature, which reduces overfitting and improves generalization.

The SVM algorithm performed slightly worse but was more interpretable due to its linear decision boundaries. Its performance was consistent across folds, highlighting its robustness.

The LSTM model, while generally used for sequential data, showed the lowest performance. This can be attributed to the limited sequential nature of the dataset and the relatively small size, which may not fully leverage LSTM's capabilities.

# Conclusion

This study highlights the strengths and limitations of different classification algorithms applied to a medical dataset. Random Forest emerged as the most effective method for this binary classification task, making it a reliable choice for similar problems. Future work could involve testing additional datasets or exploring hyperparameter tuning for further optimization.

# Code Implementation:

```python
# Import necessary libraries
import pandas as pd
import numpy as np
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split, cross_val_score, StratifiedKFold
from sklearn.metrics import confusion_matrix
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dense
from tensorflow.keras.utils import to_categorical

# Load the WDBC dataset
# File paths to the uploaded dataset
data_path = '/Users/tejasbk/Documents/1 Fall 2024/Assignments/Data mining/wdbc.data'
columns = ['ID', 'Diagnosis'] + [f'Feature_{i}' for i in range(1, 31)]  # Assuming 30 features based on typical
data = pd.read_csv(data_path, header=None, names=columns)

# Preprocess dataset
def preprocess_data(data):
    # Drop ID column as it's not useful for classification
    data = data.drop('ID', axis=1)

    # Encode the 'Diagnosis' column (B = 0, M = 1)
    data['Diagnosis'] = data['Diagnosis'].map({'B': 0, 'M': 1})

    # Separate features and target
    X = data.drop('Diagnosis', axis=1)
    y = data['Diagnosis']
    return X, y

X, y = preprocess_data(data)

# Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Initialize cross-validation
kf = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)

# Helper function to calculate metrics
def calculate_metrics(y_true, y_pred):
    tn, fp, fn, tp = confusion_matrix(y_true, y_pred).ravel()
    accuracy = (tp + tn) / (tp + tn + fp + fn)
    precision = tp / (tp + fp) if (tp + fp) > 0 else 0
    recall = tp / (tp + fn) if (tp + fn) > 0 else 0
    fpr = fp / (fp + tn) if (fp + tn) > 0 else 0
    fnr = fn / (fn + tp) if (fn + tp) > 0 else 0
    return {
        'Accuracy': accuracy,
        'Precision': precision,
        'Recall': recall,
        'FPR': fpr,
        'FNR': fnr
    }

# Random Forest Implementation
rf_model = RandomForestClassifier(random_state=42)
rf_metrics = []
for train_index, test_index in kf.split(X, y):
    X_train_fold, X_test_fold = X.iloc[train_index], X.iloc[test_index]
    y_train_fold, y_test_fold = y[train_index], y[test_index]

    rf_model.fit(X_train_fold, y_train_fold)
    y_pred_fold = rf_model.predict(X_test_fold)
    rf_metrics.append(calculate_metrics(y_test_fold, y_pred_fold))
rf_table = pd.DataFrame(rf_metrics)

# Add a gap for clarity
print("\n\n")

# SVM Implementation
svm_model = SVC(kernel='linear', random_state=42)
svm_metrics = []
for train_index, test_index in kf.split(X, y):
    X_train_fold, X_test_fold = X.iloc[train_index], X.iloc[test_index]
    y_train_fold, y_test_fold = y[train_index], y[test_index]
```

```python
    svm_model.fit(X_train_fold, y_train_fold)
    y_pred_fold = svm_model.predict(X_test_fold)
    svm_metrics.append(calculate_metrics(y_test_fold, y_pred_fold))
svm_table = pd.DataFrame(svm_metrics)

# Add a gap for clarity
print("\n\n")

# LSTM Implementation
lstm_metrics = []
X_lstm = X.values.reshape(X.shape[0], 1, X.shape[1])  # Reshape for LSTM

# Define the LSTM model once
lstm_model = Sequential([
    LSTM(64, input_shape=(X_lstm.shape[1], X_lstm.shape[2]), activation='relu'),
    Dense(1, activation='sigmoid')
])

lstm_model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])


for train_index, test_index in kf.split(X_lstm, y):

    X_train_fold, X_test_fold = X_lstm[train_index], X_lstm[test_index]
    y_train_fold, y_test_fold = y[train_index], y[test_index]

    # Train the model
    lstm_model.fit(X_train_fold, y_train_fold, epochs=3, batch_size=32, verbose=1)

    # Predict and calculate metrics
    y_pred_fold = (lstm_model.predict(X_test_fold) > 0.5).astype(int).flatten()
    lstm_metrics.append(calculate_metrics(y_test_fold, y_pred_fold))


lstm_table = pd.DataFrame(lstm_metrics)
```

```python
# Print metrics after epoch logs
print("\nRandom Forest Metrics:\n", rf_table.to_markdown(index=False))
print("\n\n")
print("SVM Metrics:\n", svm_table.to_markdown(index=False))
print("\n\n")
print("LSTM Metrics:\n", lstm_table.to_markdown(index=False))

# Add a gap for clarity
print("\n\n")

# Collect results
rf_avg_metrics = rf_table.mean().to_dict()
svm_avg_metrics = svm_table.mean().to_dict()
lstm_avg_metrics = lstm_table.mean().to_dict()

results = pd.DataFrame({
    'Random Forest': rf_avg_metrics,
    'SVM': svm_avg_metrics,
    'LSTM': lstm_avg_metrics
})

print("Comparison of Classification Algorithms:\n", results.to_markdown())
```

# Output:

```
Epoch 1/3
16/16 [==============================] - 0s 735us/step - loss: 28.1423 - accuracy: 0.6289
Epoch 2/3
16/16 [==============================] - 0s 666us/step - loss: 3.8112 - accuracy: 0.4102
Epoch 3/3
16/16 [==============================] - 0s 613us/step - loss: 0.7630 - accuracy: 0.7617
2/2 [==============================] - 0s 776us/step
Epoch 1/3
16/16 [==============================] - 0s 623us/step - loss: 0.3429 - accuracy: 0.8750
Epoch 2/3
16/16 [==============================] - 0s 587us/step - loss: 0.2900 - accuracy: 0.9043
Epoch 3/3
16/16 [==============================] - 0s 556us/step - loss: 0.2874 - accuracy: 0.9102
2/2 [==============================] - 0s 731us/step
Epoch 1/3
16/16 [==============================] - 0s 552us/step - loss: 0.2944 - accuracy: 0.9004
```

Random Forest Metrics:

| Accuracy | Precision | Recall | FPR | FNR |
|---------:|----------:|--------:|---------:|---------:|
| 0.947368 | 0.913043 | 0.954545 | 0.0571429 | 0.0454545 |
| 0.982456 | 0.956522 | 1 | 0.0285714 | 0 |
| 0.964912 | 0.952381 | 0.952381 | 0.0277778 | 0.047619 |
| 0.912281 | 0.944444 | 0.809524 | 0.0277778 | 0.190476 |
| 0.947368 | 0.95 | 0.904762 | 0.0277778 | 0.0952381 |
| 0.964912 | 1 | 0.904762 | 0 | 0.0952381 |
| 0.964912 | 0.952381 | 0.952381 | 0.0277778 | 0.047619 |
| 0.947368 | 0.875 | 1 | 0.0833333 | 0 |
| 0.929825 | 0.947368 | 0.857143 | 0.0277778 | 0.142857 |
| 1 | 1 | 1 | 0 | 0 |

SVM Metrics:

| Accuracy | Precision | Recall | FPR | FNR |
|---------:|----------:|--------:|---------:|---------:|
| 0.964912 | 0.954545 | 0.954545 | 0.0285714 | 0.0454545 |
| 0.947368 | 0.88 | 1 | 0.0857143 | 0 |
| 0.947368 | 0.95 | 0.904762 | 0.0277778 | 0.0952381 |
| 0.912281 | 0.944444 | 0.809524 | 0.0277778 | 0.190476 |
| 0.964912 | 1 | 0.904762 | 0 | 0.0952381 |
| 0.947368 | 0.95 | 0.904762 | 0.0277778 | 0.0952381 |
| 0.964912 | 0.952381 | 0.952381 | 0.0277778 | 0.047619 |
| 0.947368 | 0.909091 | 0.952381 | 0.0555556 | 0.047619 |
| 0.964912 | 1 | 0.904762 | 0 | 0.0952381 |
| 0.982143 | 1 | 0.952381 | 0 | 0.047619 |

LSTM Metrics:

| Accuracy | Precision | Recall | FPR | FNR |
|---------:|----------:|--------:|---------:|---------:|
| 0.859649 | 0.888889 | 0.727273 | 0.0571429 | 0.272727 |
| 0.877193 | 0.941176 | 0.727273 | 0.0285714 | 0.272727 |
| 0.964912 | 0.952381 | 0.952381 | 0.0277778 | 0.047619 |
| 0.894737 | 0.941176 | 0.761905 | 0.0277778 | 0.238095 |
| 0.824561 | 0.703704 | 0.904762 | 0.222222 | 0.0952381 |
| 0.859649 | 1 | 0.619048 | 0 | 0.380952 |
| 0.877193 | 0.9375 | 0.714286 | 0.0277778 | 0.285714 |
| 1 | 1 | 1 | 0 | 0 |
| 0.947368 | 1 | 0.857143 | 0 | 0.142857 |
| 0.982143 | 0.954545 | 1 | 0.0285714 | 0 |

Comparison of Classification Algorithms:

| | | Random Forest | SVM | LSTM |
|:----------|---|----------------:|---------:|---------:|
| Accuracy | | 0.95614 | 0.954355 | 0.908741 |
| Precision | | 0.949114 | 0.954046 | 0.931937 |
| Recall | | 0.93355 | 0.924026 | 0.826407 |
| FPR | | 0.0307937 | 0.0280952 | 0.0419841 |
| FNR | | 0.0664502 | 0.075974 | 0.173593 |