# Project Assignment

# Task C2 – Report

**Rajapaksha Mudiyanselage Udara Ishan Bandara Embawa**
**105107155**

## Executive Summary

This report documents the enhancement of the stock prediction system from version 0.1 to version 0.2, addressing critical limitations in data processing and feature utilization. The primary improvement involves implementing multi-feature processing capabilities, allowing the LSTM model to utilize Open, High, Low, Close, Adjusted Close, and Volume data instead of only Close prices.

The original code's main limitations were:

- Single feature utilization (Close price only)

- Manual date selection requirements

- Limited data preprocessing capabilities

The enhanced version addresses these issues through improved data processing functions and multi-feature LSTM implementation.

## Key Code Modifications

Multi-Feature Data Scaling Function

Original Implementation:

python

```python
def scale_close_series(df):
    data = df[["Close"]].values
    scaler = MinMaxScaler(feature_range=(0, 1))
    scaled = scaler.fit_transform(data)
    return scaled, scaler
```

Enhanced Implementation:

python

```python
def scale_features(df, scaler_type="minmax"):
    # candidate features
    feature_candidates = ["Open", "High", "Low", "Close", "Adj Close", "Volume"]
    # keep only those actually in df
    feature_cols = [c for c in feature_candidates if c in df.columns]

    X = df[feature_cols].values
    y = df["Close"].values  # target is still Close

    if scaler_type == "minmax":
        scaler = MinMaxScaler()
    else:
        scaler = StandardScaler()

    X_scaled = scaler.fit_transform(X)

    return X_scaled, y, scaler, feature_cols
```

Key Improvements:

- Dynamic Feature Detection: The list comprehension [c for c in feature_candidates if c in df.columns] automatically detects available features, preventing crashes when datasets have different column structures.

- Flexible Scaling Options: The scaler_type parameter allows choosing between MinMaxScaler (0-1 normalization) and StandardScaler (z-score normalization).

- Enhanced Return Values: Returns the feature column list and scaler object for proper inverse transformation later.

Multi-Feature Sequence Generation

      Enhanced Sequence Function:

python

```python
def make_sequences_multifeature(X, y, time_step=60):
    Xs, ys = [], []
    for i in range(time_step, len(X)):
        Xs.append(X[i-time_step:i, :])   # use all features
        ys.append(y[i])                  # target is Close
    return np.array(Xs), np.array(ys)
```

Technical Explanation:

- Multi-Dimensional Slicing: X[i-time_step:i, :] captures all features (: in second dimension) for each time window, creating sequences with shape (time_steps, n_features).

- Target Preservation: While input uses multiple features, the target remains Close price for consistent prediction objectives.

Complex Inverse Transformation

      Most Technically Challenging Section:

python

```python
dummy = np.zeros((len(pred_scaled), len(feature_cols)))
dummy[:, feature_cols.index("Close")] = pred_scaled[:, 0]
pred = scaler.inverse_transform(dummy)[:, feature_cols.index("Close")]
```

Research-Required Explanation:

- Scaler Constraint: Scikit-learn scalers require the same number of features for inverse transformation as used during fitting. Since the scaler was fitted on all features, inverse transformation needs a complete feature array.

- Dummy Array Strategy: np.zeros() creates a matrix with correct dimensions, filling non-target features with zeros.

- Selective Population: feature_cols.index("Close") dynamically locates the Close column position, ensuring code works regardless of feature order.

- Two-Step Extraction: After inverse transformation, only the Close column values are extracted, converting back to actual price scale.

# Research-Intensive Components

LSTM Input Shape Requirements

Research Finding: LSTM layers in Keras require 3D input tensors with shape (batch_size, time_steps, features). The original code used (batch_size, 60, 1), while the enhanced version uses (batch_size, 60, n_features).

Implementation Impact:

python

X_train = X_train.reshape((X_train.shape[0], X_train.shape[1], 1))  *# Original*

*# vs*

*# No reshaping needed - already correct shape from make_sequences_multifeature*

Feature Engineering in Financial Markets

Research Insight: Financial markets exhibit complex relationships between different price indicators:

- Volume-Price Relationship: High volume often precedes significant price movements

- Price Range Patterns: High-Low spreads indicate volatility and market sentiment

- Opening Gap Analysis: Differences between previous close and next open reveal overnight sentiment

Normalization Strategy Selection

Research-Based Decision: MinMaxScaler chosen over StandardScaler for financial data because:

- Financial data often has non-normal distributions

- LSTM networks perform better with bounded inputs (0-1 range)

- Preserves relative relationships between different price levels

## Implementation Challenges and Solutions

Tensor Dimension Management

Challenge: Managing 3D tensor operations for multi-feature sequences while maintaining prediction accuracy.

Solution: Systematic approach using .shape attributes and dynamic indexing to ensure dimensional consistency throughout the pipeline.

Future Prediction with Unknown Features

Challenge: Predicting future prices when other feature values (Open, High, Low, Volume) are unknown.

Solution Implemented:

python

```python
# create dummy row with n_features
dummy_features = np.zeros((1, 1, X_scaled.shape[1]))
dummy_features[0, 0, close_idx] = next_scaled

# slide the window
cur = np.append(cur[:, 1:, :], dummy_features, axis=1)
```

This creates a sliding window where only the predicted Close price is known, while other features are zero-padded.

Scaler State Management

Challenge: Maintaining scaler objects for proper inverse transformation across different feature sets.

Solution: Returning scaler and feature column information as function outputs, enabling proper inverse transformation later in the pipeline.

# Performance and Accuracy Implications

Model Capacity Enhancement

The multi-feature approach increases model learning capacity by:

- Providing richer input context through multiple market indicators

- Enabling the LSTM to learn cross-feature correlations

- Potentially improving prediction accuracy through comprehensive market representation

Computational Considerations

Trade-offs:

- Increased Training Time: Proportional to number of features used

- Higher Memory Requirements: 3D tensors consume more memory than 2D

- Enhanced Learning Potential: More features can lead to better pattern recognition

# Code Quality Improvements

Modularity and Maintainability

The enhanced code demonstrates improved software engineering practices:

- Function Separation: Clear distinction between data processing and model training

- Parameter Flexibility: Configurable scaling methods and feature selection

- Error Prevention: Robust handling of missing or unexpected data columns

Scalability Features

- Dynamic Feature Handling: Code adapts to different datasets automatically

- Configurable Architecture: Easy modification of LSTM parameters and structure

- Extensible Design: Simple addition of new features or processing methods

## Technical Learning Outcomes

Advanced NumPy Operations

Gained proficiency in:

- Multi-dimensional array manipulation and slicing

- Tensor reshaping for deep learning applications

- Dynamic indexing and array construction

Deep Learning Pipeline Development

Enhanced understanding of:

- LSTM input requirements and tensor shapes

- Feature scaling impact on neural network performance

- Time series prediction methodology for financial data

Financial Data Processing

Developed expertise in:

- Multi-feature financial dataset handling

- Proper normalization techniques for market data

- Time series forecasting with multiple indicators

## Conclusion

The transition from single-feature to multi-feature processing represents a significant advancement in the stock prediction system's capability. The implementation successfully addresses the original limitations while maintaining code clarity and functionality.

Key Achievements:

- Successfully implemented multi-feature LSTM processing

- Solved complex inverse transformation challenges

- Created robust, scalable data processing pipeline

- Maintained prediction accuracy while enhancing model input richness

Most Complex Technical Component: The inverse transformation process required extensive research into scikit-learn scaler mechanics and NumPy array manipulation, resulting in an elegant dummy array solution that maintains dimensional consistency while extracting relevant predictions.

The enhanced code provides a solid foundation for advanced financial prediction models and demonstrates mastery of both machine learning principles and practical implementation challenges in time series forecasting.