**Name: Deepak Kumar B**

**Date: 26/05/2024**

# Time Series Prediction Using Deep Learning

## Introduction

The primary objective of this project is to create a predictive model using deep learning techniques to forecast future sales based on historical time series data. The dataset comprises hourly sales data from a major retail chain spanning from 2015 to 2020. Accurate sales forecasting is critical for inventory management, financial planning, and optimizing operations in retail.

## Data Exploration and Preprocessing

### Dataset Description

The dataset includes the following columns:

- **Date:** This column records the timestamp of each sales transaction.

- **Sales:** This column represents the total sales value in USD.

### Loading the Data

The initial step involved loading the dataset into a pandas, DataFrame. This is a crucial step as it allows us to inspect the first few rows of the data to understand its structure, identify any inconsistencies, and get an overview of the data types in each column.

## Data Cleaning and Preparation

### Cleaning Column Names

It was essential to ensure that the column names were free from any leading or trailing whitespace. This was achieved by stripping the column names. Clean column names are vital for consistency and to avoid errors during further data manipulation and analysis.

## Converting Date to Datetime

The 'Date' column was converted to a datetime format. This transformation is critical as it allows for efficient handling of time series data. The 'Date' column was then set as the index of the DataFrame, facilitating time-based operations and analyses.

## Handling Missing Data

Handling missing data is a significant step in data preprocessing. Any rows containing missing values were removed to ensure the dataset was clean. This step helps in maintaining the integrity of the data, ensuring that the model is trained on complete and accurate records.

## Data Visualization

Visualization of the sales data was performed to identify trends, seasonal patterns, and any anomalies. A line plot of the sales data over time revealed the overall trend, periodic fluctuations, and potential seasonal effects. This visualization provided valuable insights into the data's behavior, which is essential for selecting appropriate modeling techniques.

## Scaling the Sales Values

Normalization of the sales values was performed using MinMaxScaler, scaling the data to a range between 0 and 1. This step is crucial because it ensures that the data is within a standard range, improving the performance and convergence speed of the neural network. Deep learning models, particularly those involving gradient descent, benefit from scaled inputs as they stabilize and speed up the learning process.

# Model Selection and Implementation

## Model Choice

For this time series forecasting task, Long Short-Term Memory (LSTM) networks were chosen. LSTMs are a type of recurrent neural network (RNN) designed to capture long-term

dependencies in sequential data. They are particularly effective in handling the vanishing gradient problem, making them suitable for time series data with its inherent temporal dependencies.

## Creating the Dataset for LSTM

The dataset was transformed to create input-output pairs suitable for LSTM models. This involved defining a 'look-back' period, which specifies the number of previous time steps used to predict the next time step. This transformation is critical as it structures the data into sequences that the LSTM can learn from.

## Splitting the Data

The data was split into training and testing sets to evaluate the model's performance on unseen data. Typically, 80% of the data was used for training, and the remaining 20% was reserved for testing. This split ensures that the model is not only fitted to the historical data but is also evaluated on its ability to generalize to new, unseen data.

## Model Architecture

The LSTM model architecture was carefully designed:

- **Input Layer:** This feeds the input data into the model.

- **LSTM Layers:** The model includes two LSTM layers, each with 50 units. The first LSTM layer is configured to return sequences, allowing the stacking of another LSTM layer. The second LSTM layer processes these sequences and extracts relevant features for prediction.

- **Dense Layers:** These layers transform the LSTM outputs to the desired shape. They act as fully connected layers that combine the features learned by the LSTM layers.

- **Compilation:** The model was compiled using the Adam optimizer and the mean squared error (MSE) loss function. The Adam optimizer is chosen for its efficiency and effectiveness in training deep learning models, and MSE is suitable for regression tasks as it measures the average squared differences between predicted and actual values.

## Training the Model

The model was trained using the training dataset. Key parameters such as batch size and number of epochs were specified. The batch size determines how many samples are processed before the model's internal parameters are updated, while the number of epochs

defines how many times the learning algorithm will work through the entire training dataset. The training process was monitored by recording the training and validation loss, providing insights into how well the model is learning and if it is overfitting or underfitting.

# Evaluation Metrics

## Predictions and Inverse Transform

Once the model was trained, predictions were made for both the training and testing sets. These predicted values were then scaled back to their original range using the inverse transform of the MinMaxScaler. This step is crucial for interpreting the results in their original context and for calculating meaningful performance metrics.

## Performance Metrics

The model's performance was evaluated using several metrics:

- **Root Mean Square Error (RMSE):** This metric measures the square root of the average squared differences between predicted and actual values. It provides a measure of the magnitude of errors in the predictions.

- **Mean Absolute Error (MAE):** This metric measures the average absolute differences between predicted and actual values. It gives an idea of how close the predictions are to the actual values on average.

- **$R^2$ Score:** This score indicates the proportion of the variance in the dependent variable that is predictable from the independent variables. It provides a measure of how well the model's predictions match the actual data.

These metrics collectively offer a comprehensive assessment of the model's accuracy and reliability.

# Visualizing Predictions

Predictions were visualized alongside actual sales values to provide a clear picture of the model's performance. Separate plots for the training and testing periods helped in understanding how well the model learned from the training data and how accurately it predicted the testing data. Visualization is a powerful tool for highlighting the model's strengths and identifying areas where it may need improvement.

# Conclusion

The project successfully developed an LSTM-based deep learning model for time series forecasting of retail sales. The model demonstrated the ability to capture temporal dependencies in the data and make accurate predictions. Key steps included thorough data preprocessing, careful model selection and design, and rigorous evaluation using various metrics.

This process underscores the importance of data preparation, appropriate model selection, and thorough evaluation in developing effective predictive models. The insights gained from this project can be applied to similar time series forecasting tasks in different domains, illustrating the versatility and power of deep learning techniques in predictive analytics.