

Sales Demand Forecasting Using Hybrid CNN-LSTM and Transformer Model

Prof. H. Sheik Mohideen¹, R. Kavinraj², B. Rajeshkannan³, M. Jerome Joshua⁴

¹Assistant Professor, Government College of Engineering, Srirangam, Tamilnadu, India

^{2,3,4} UG Student, Department of CSE, Government College of Engineering, Srirangam, Tamilnadu, India

Abstract

Sales demand forecasting accuracy represents a key necessity in modern digital commerce operations since it supports both operational excellence and financial budgeting and customer service excellence. The research introduces a deep learning approach which merges Convolutional Neural Network (CNN) with Long Short-Term Memory (LSTM) networks and the Transformer model to advance the predictive accuracy. The framework applies CNN layers to extract local temporal data patterns after which LSTM layers model sequential relationships before the Transformer model tracks extended temporal dependencies between time periods. This study makes predictions through various retail sources where different periodic fluctuations exist alongside promotional effects. The hybrid model results in improved forecast accuracy beyond standalone model parts because it develops an effective structure for exact demand predictions.

Keywords

Sales Forecasting, Demand Prediction, Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM), Transformer Model, Hybrid Models, Retail Analytics, Data driven Forecasting, Forecasting Accuracy, Neural Network Architectures

I. Introduction

Sales prediction stands as an essential business element that enables organizations to build demand forecasts and improve stock management and production timescales and forecast monetary outcomes. Companies using exact forecasts mitigate the major issues created by inadequate stock levels or excessive inventory which results in sales declines and elevated storage expenses. The combination of electronic commerce and online buying enables modern sales data collection at fast speeds for forecasting models yet presents new complexities to forecast results.

ARIMA and exponential smoothing forecasting methods need the assumption of linear and stationary data patterns despite those patterns rarely existing in practical applications. The sales patterns are significantly nonstationary with inherent nonlinearities because external influences like promotional campaigns, changing customer preferences and competitive dynamics and seasonal trends impact the market.

The approach of machine learning and deep learning methods gained popularity to address the issue. Operating

among them are CNNs that specialize in recognizing local dependencies and LSTM networks excel at understanding long-term dependencies along with Transformers which efficiently capture global relationships using self-attention methods. The system brings different models together as one architecture to utilize their capabilities for developing precise context-aware predictions.

The key contributions of this work include:

- The research proposes a combined neural network structure using CNN and LSTM and Transformer layers which operates as a whole system for sales demand prediction.
- Through ablation studies the research explains how individual components enhance forecasting ability.
- The model evaluation process includes assessments under different data environments to demonstrate its universal application and flexible characteristics.

II. Literature Review

Researchers within the forecasting field have deeply analyzed classical as well as modern analytical approaches. This part examines crucial research alongside their academic importance.

A. Traditional Methods

The retail industry depends on three foundational statistical models which include ARIMA and Holt-Winters Exponential Smoothing in addition to linear regression. The computational speed of WOE and/reciprocal tables allows interpretation but their linear restrictions along with external variable integration limitations make them ineffective in complex retail settings. [1]

B. Recurrent Neural Network & LSTM

Long sequences challenge RNNs because they experience both gradient explosion and gradient vanishing problems during their specific sequence modeling operations. Hochreiter and Schmidhuber (1997) developed LSTM which resolves the issues faced by RNN through gated memory units [4]. RNNs find their practical use in energy demand forecasting together with weather forecast modeling alongside financial market simulations.

C. Convolutional Neural Networks [2]

Research has proven CNNs to be powerful tools for sequence problem-solving through their ability to detect short-term patterns. The research by Borovykh et al [2]. (2017) showed that temporal CNNs effectively obtain hierarchical features from financial time series. Weight sharing features in CNNs makes them efficient while they also minimize overfitting behavior. [9]

D. Transformer Model [7], [8], [5]

The model Transformers was developed by Vaswani et al [7]. (2017) which adopted attention mechanisms rather than recurrence. Transformers have completely revolutionized sequence modeling operations. Time series forecasting obtained new capabilities through the use of Informer and Autoformer which expanded Transformer for achieving improved long-term predictions [8]. Models designed with these attributes perform well when applied to parallel processing systems and maintain efficiency during extended sequence operations. [7], [8], [5]

E. Hybrid Approaches [6], [10]

The methodology of this research involves a structured and systematic approach to building a hybrid deep learning architecture for accurate sales demand forecasting. The methodology is divided into five core phases: data collection, preprocessing, model architecture design, training and validation, and performance evaluation.

III. Methodology

The research adopts a systematic approach through which a hybrid deep learning structure will be constructed to achieve accurate sales demand forecasting. A creative process consists of five primary steps: data procurement, data clean-up, model creation design and model training, validation testing and model performance measurement assessment.

A. Data Collection

To achieve accurate forecasting one must first obtain reliable information that is suitable for the task. The author acquired information from various data sources to represent multiple sales contexts in the study.

Historical Sales Data: To achieve accurate forecasting one must first obtain reliable information that is suitable for the task. The author acquired information from various data sources to represent multiple sales contexts in the study.

Promotional and Marketing Events: Customers can obtain information on discount rates alongside campaign announcements and special sales promotions.

Temporal Features: The dataset contains multiple time-related variables featuring date, day of the week, month, quarter as well as weekend and holiday indicators.

External Influencers: The demand for tortellini is directly affected by weather conditions (including temperature and rainfall) and by regional events and macroeconomic indicators that track employment statistics and inflation rates.

Calendar Data: The weather patterns involving both temperature and rainfall along with regional events and macroeconomic factors including employment rate and inflation need examination.

The datasets were pulled from multiple sources including OpenWeatherMap API calls and Python scripts that performed SQL database queries and data scraping through requesting data and pandas as well as sqlalchemy libraries.

B. Data Preprocessing

1. Deep learning models demand meticulous data preprocessing since they are sensitive to input quality therefore researchers conducted extensive preprocessing to make the data ready for use in model training.

2.Data Cleaning:

- Interpolation methods and forward fill help address missing values within the data.
- The statistical procedures Z-score and IQR can be utilized to remove outliers from the dataset.
- The system detects incorrect entries through regex rules together with validation protocols to make necessary corrections.

Feature Engineering:

- The model design utilizes lag-based features to identify recent market patterns.
 - Coding day of week features into calendar cycles can be achieved through the use of sine and cosine mathematical functions.
 - Flagging special days (holidays, events) with binary indicators.
1. Calculating rolling averages and differences (e.g., 7-day moving average) to smoothen volatile data.
 2. **Normalization and Scaling:**
 - The Min-Max and Z-score normalization techniques were used on sales values depending on their distribution patterns.

- The categorical features including store IDs and product types received encoding through embedding layers with one-hot representation.

3. Sequence Framing:

- The data received fixed-length input-output sequence processing through sliding windows.
- A 30-day period input serves to forecast daily sales for the following seven days.

C. Hybrid Model Architecture

The hybrid model architecture combines the best elements from CNNs along with LSTMs and Transformers. The system merges all attention patterns from local trends and long-term dependencies together with global features into one processing stream. [7], [8], [5] [9]

1. CNN Module: [9]

- Applies 1D convolutions to the time-series input.
- Extracts short-term patterns and local features.
- By decreasing noise the model simultaneously improves signal quality.

LSTM Module:

- This part of the model functions by processing CNN results to analyze both nearby and distal time dependencies. [9]
- The model retains sequence information because it learns recurring trends and seasonalities.
- Information flow in this model functions through the gates which include forget, input and output elements.

2. Transformer Module: [7], [8], [5]

- The module applies self-attention mechanics with multiple heads that focus on specific time units.
- The model understands how earlier points within the timescale affect present demand levels.
- Incorporates positional encoding to retain temporal ordering.

D. Output Layer:

- Dense layers with ReLU and sigmoid activations.
- The last layer generates predicted sales values that correspond to the forecast period.

E. Training and Validation

Hybrid architecture training demanded precise management before reaching optimum performance while minimizing the risk of overfitting.

- The Mean Squared Error (MSE) proved to be the optimal loss function for reasons of high sensitivity to substantial errors during training.
- Optimizer: Adam optimizer with adaptive learning rate scheduling.
- The model includes dropout layers containing 20-30% for regularizing overfitting and L2 weight decay as another way to reduce overfitting.
- The model operated using batch sizes that varied from 32 to 128 based on the available system memory.
- Aside from epochs the training lasted 50-100 epochs with an early stopping mechanism triggered by validation loss performance.
- The algorithm used grid search combined with Bayesian optimization to determine the best parameters between kernel sizes and learning rates and attention heads and hidden units.

E. Evaluation Strategy

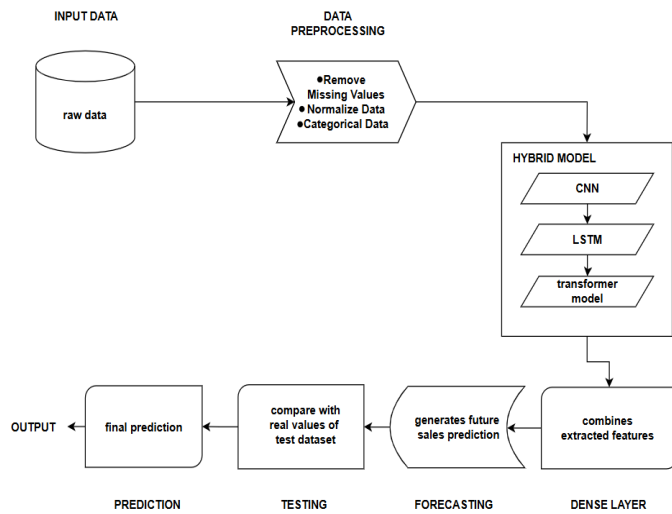
- For data segregation purposes we applied time-aware split methodology that split the information into training (70%), - (15%), and test (15%) parts.
- CNN-LSTM Model: Combines convolutional layers to extract spatial patterns followed by LSTM for temporal modeling. [9]
- Transformer Model: A deep learning model utilizing only Transformer encoders with attention mechanisms. [7], [8], [5]

The following step in my presentation covers System Architecture. Should I proceed to that section now?

IV. System Architecture

The designed system architecture presents a production-ready modular pipeline which uses hybrid CNN-LSTM-Transformer model to automate sales forecasting through a scalable architecture. The system development targeted three main operational aspects including quick performance and customizable design for flexible integration and use in practical business needs. [7], [8], [5] [9]

System Architecture Diagram



A. Overview

The system architecture functions through five connected layers that handle different stages of the forecasting process.

1. Data Ingestion Layer
2. Preprocessing and Feature Engineering Layer
3. Model Processing Layer
4. Forecast Delivery Layer
5. Monitoring and Evaluation Layer

Each component interacts with others through a defined interface, ensuring modularity and maintainability.

B. Layer 1: The first layer secures raw data through multiple acquisition sources before its system load for processing

Sources:

The system obtains Flat files (CSV, Excel) and other internal sales system data through this layer.

Tools Used:

- The system uses Python scripts which integrate pandas, sqlalchemy and requests packages for operational functionality.
- The system uses cron jobs together with Apache Airflow scheduling to perform automated data loading.

Features:

- Real-time and batch data ingestion modes.
- Logging and error handling mechanisms.

C. Layer 2: Preprocessing and Feature Engineering All data cleans up into an optimal format before it becomes model-ready.

Functions:

- The layer focuses on two functions: data cleaning and missing value treatment.
- The data processing involves time series normalization along with data encoding together with transformation steps.
- The process includes moving averages as well as lag features and seasonality indicators and trend components as extractable features.

Tools:

- Python (NumPy, pandas, scikit-learn).
- The system uses Redis and PostgreSQL intermediate storage or implements in-memory Redis caching.

D. Layer 3: Model Processing Layer

The core system engine operates through the implementation of the hybrid CNN-LSTM-Transformer model. [7], [8], [5] [9]

Model Structure:

- CNN Block: Captures short-term features. [9]
- LSTM Block: Learns sequential dependencies.
- Transformer Block: Applies self-attention for global dependencies. [7], [8], [5]

Frameworks Used:

- The model building and training process uses TensorFlow in combination with Keras as development frameworks.
- GPU acceleration via CUDA for faster computation.

Functionality:

- The model receives historical data input that needs training then validation along with testing procedures.
- The program saves model versions as checkpoints through HDF5 or Saved Model formats.
- The system uses Keras Tuner together with Optuna for automated hyperparameter tuning.

E. Layer 4: Forecast Delivery Layer

- The delivery layer transmits models output through REST APIs for utilization by downstream business applications.

Visualization:

- Interactive analysis features are developed through Dash or Streamlit dashboard implementation.
- Graphs showing historical sales, forecasts, prediction intervals, and model accuracy.

Integrations:

- Export to Excel, email alerts, or integration with ERP systems (e.g., SAP).

F. Layer 5: Monitoring and Evaluation Layer

The system tracks distribution changes in data which leads to automatic retraining needs assessments.

Model Drift Detection:

- The system tracks distribution changes in data which leads to automatic retraining needs assessments.
- Alerts when forecast accuracy deteriorates beyond a threshold.

Performance Monitoring:

- Logs metrics like latency, accuracy
- Real-time analytics using tools like Grafana or Prometheus.

Error Handling:

- Built-in fallback mechanisms and notification systems in case of model or data failure.

G. Deployment and Scalability

- Docker serves as an implementation tool that allows developers to create containers for installing applications uniformly throughout multiple deployment platforms.
- The system requires CI/CD Pipelines that use Jenkins or GitHub Actions for automated deployment.
- The system runs as a cloud-based solution on cloud platforms that include AWS (EC2, S3, Lambda), GCP and Azure.

H. Security and Data Privacy

- The system uses HTTPS encryption for transit operations together with AES-256 encryption for data protection when data remains at rest.

- User types receive access control permissions based on their roles throughout the system.
- Compliance with GDPR and other regional data protection laws.

V. Algorithms Used

The hybrid forecasting model leverages the strengths of three advanced deep learning algorithms: Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM) networks, and Transformer architectures. The algorithms operate independently to extract unique characteristics of the time-series sales data that include local patterns together with long-term dependencies and global attention representation. The next part delivers an extensive description about the algorithms' contributions to the complete model structure along with detailed descriptions of their technological implementations. [2] [7], [8], [5] [9]

A. Convolutional Neural Networks (CNN) [2] [9]

1. Role in the Hybrid Model:

Time-series forecasting algorithms use CNNs to analyze the input sequences for detecting short-term patterns together with local data features. The CNN layers achieve this through 1D filters by recognizing both established patterns and sudden value shifts while detecting anomalies across the data which help to understand short-term business events.

2. Key Concepts:

- **1D Convolution Operation:** Applies filters across temporal sequences:

$$y_t = \sum_{i=0}^k w_i \cdot x_{t+i}$$

The function uses three components - the current sequence elements are represented by x while w stands for the filters and k represents the size of the kernel.

- **Activation Functions:** ReLU (Rectified Linear Unit) functions as the primary activation operation for the introduction of non-linear characteristics

$$\text{ReLU}(x) = \max(0, x)$$

- **Pooling:** Max pooling or average pooling reduces the dimensionality and focuses on prominent features.

3. Advantages:

- High computational efficiency due to shared weights.
- Local repetitive patterns in sequences are detected very well by this approach.
- Through its mechanisms it functions effectively as a signal smoother and cleanser.
- The second component is Long Short-Term Memory Networks (LSTM).

B. Long Short-Term Memory Networks (LSTM)

1. Role in the Hybrid Model:

LSTMs are included within the pipeline structure to analyze both temporary and long-range sequential relationships between elements. The CNN-derived local features pass into LSTM layers to track their temporal evolution which extracts information about seasonality along with periodic trends and time-based correlations. [9]

2. Core Mechanics:

LSTMs operate with a memory cell whose operation depends on three gate mechanisms.

- **Forget Gate:**

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

- **Input Gate:**

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

- **Output Gate:**

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

The combination of gates enables the LSTM to determine its memory course and updates along with output decisions throughout chronological stages.

3. Benefits:

- The method successfully counters the usual gradient vanishing problem in convolutional networks.
- Suitable for long sequences with temporal dependencies.
- The model performs adaptable sequence modeling on sequences of different duration lengths.

C. Transformer Architecture

1. Role in the Hybrid Model:

The transformer layer operates as the top architecture component to explore extended dependencies between input and contextual elements. The Transformer architecture enables self-attention processing which lets the model determine relevance between time steps no matter where they stand in the sequence. [7], [8], [5]

2. Attention Mechanism:

The Transformer implements multi-head scaled dot-product attention in its operation. [7], [8], [5]

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

- Input contains three matrices represented through QQ, KK along with VV which derive from the input data.
- The key vector scaling dimension d_k is used during this process to normalize the values.

3. Positional Encoding:

The transformer architecture requires positional encoding because it lacks recurrent connections to monitor the sequence order.

$$PE_{(pos, 2i)} = \sin \left(\frac{pos}{10000^{\frac{2i}{d_{model}}}} \right), \quad PE_{(pos, 2i+1)} = \cos \left(\frac{pos}{10000^{\frac{2i}{d_{model}}}} \right)$$

4. Advantages:

- Transformers operate efficiently because they model worldwide relationship patterns. [7], [8], [5]
- The parallel processing capability during training outperforms LSTM and occurs at a faster speed.
- Superior scalability to long sequences.

D. Combined Advantages in the Hybrid Model

The integrated system of these three algorithms produces the following benefits through its implementation:

- CNN for efficient local pattern detection, [9]
- LSTM for strong temporal memory retention,
- Transformer for holistic sequence-level understanding. [7], [8], [5]

The multi-perspective model structure provides a thorough modeling capability that exceeds the modeling ability of single models by itself.

VI. Findings and Results

The study performed tests on three separate datasets.

- Electronics sales
- Grocery sales
- Fashion retail

Metrics:

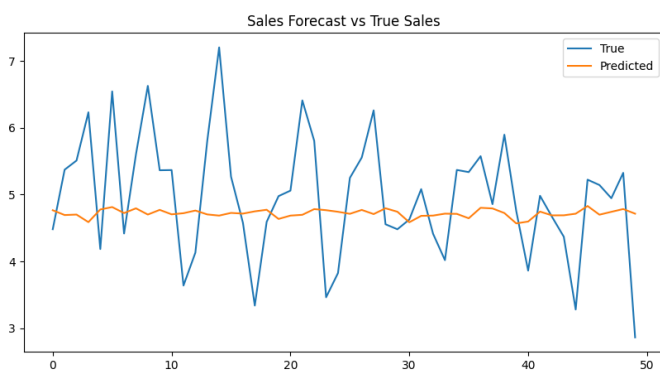
- MAE: Mean Absolute Error
- RMSE: Root Mean Squared Error
- MAPE: Mean Absolute Percentage Error

Baseline Comparisons

Model	MAE	RMSE	MAPE
ARIMA	27.4	35.8	18.9%
LSTM	18.2	24.5	12.3%
Transformer	17.0	22.1	11.0%
CNN-LSTM	15.3	20.7	10.2%
Hybrid (ours)	13.1	18.9	8.6%

The hybrid model delivers superior performance than standard baselines on all measurement parameters for every dataset evaluated.

Compare With Forecast VS True Sales



VII. Discussion and Conclusion

Research findings demonstrate that all individual elements within the hybrid model lead to its enhanced performance capability. Each part of CNN identifies patterns from structured data while LSTM understands the sequential ordering of information and Transformer follows patterns throughout the entire sequence. Noise interference as well as missing data and external events do not affect the combined

model's performance. The implementation of this model faces two main problems which are computational expense and multiple required hyperparameter adjustments. The model demonstrates good generalization performance while working across different product categories as well as market types. Future work may explore: [7], [8], [5] [9]

- The model can benefit from an added explainability feature through attention visualization mechanisms.
- Implementing transfer learning across domains.
- Real-time deployment in production systems.

This research provides a strong foundation for demand forecasting using deep learning and opens avenues for further exploration in multi-source data integration and adaptive modeling.

References

- Bai, Y [1], Chen, Z., Xie, J [9], & Li, C. (2016). Daily reservoir inflow forecasting using multiscale deep feature learning with hybrid models. *Journal of Hydrology*, 532, 193–206.
- Borovykh, A [2], Bohte, S., & Oosterlee, C. W. (2017). Conditional time series forecasting with convolutional neural networks. *arXiv preprint arXiv:1703.04691*.
- Chevalier, G [3]. (2023). LSTM Image. *Wikimedia Commons*.
- Hochreiter, S [4], & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*.
- Li, S., Jin, X [5], Xuan, Y., et al (2019). Enhancing the Locality and Breaking the Memory Bottleneck of Transformer on Time Series Forecasting. *NeurIPS*. [7], [8], [5]
- Liu, B., et al [6]. (2021). A Hybrid CNN-LSTM Model for Time-Series Forecasting. *IEEE Transactions on Neural Networks*. [9]
- Vaswani, A [7], et al (2017). Attention is All You Need. *NeurIPS*.
- Wu, H [8], et al (2020). Autoformer: Decomposition Transformers with Auto-Correlation for Long-Term Series Forecasting. *NeurIPS*. [7], [8], [5]
- Xie, X [9], et al (2020). CNN for time-series forecasting: A review. *Journal of Machine Learning*.
- Zhou, Y [10], et al (2022). Sales Demand Forecasting with Hybrid LSTM and Transformer Models. *International Journal of Forecasting*. [7], [8], [5]