

PROJECT REPORT

Deep Learning for Demand Forecasting in Retail Using Hybrid CNN-RNN and Transformer Models

NAME: Raghulchellapandian Senthil Kumaran.

NAME: Dongyoон Shin

UBID: raghulch

UBID: dongyoон

Abstract:

Precise time series predictions are essential in the retail industry, where inventory, supply chain and strategic planning depend on accurate forecasting. In this task, we study the capability of deep learning models in predicting future sales at multiple time steps ahead in the future. We test architectures (LSTM, GRU, Transformer, FFNN, hybrid CNN-RNN and basic Linear model) on multivariate time series data. The models are evaluated not just on clean data, but also in a noisy setting, to mimic the real-world uncertainty and quantify the resistance. Using a well-defined experimental configuration and evaluation criteria, we explore model strengths and limitations in two kinds of scenes, ideal and noisy.

1. PROBLEM STATEMENT:

In the fast changing and competitive environment of retail industry, predicting future demand is a significant aspect to ensure right stock level, to minimize operation cost and to enhance customer satisfaction. Classical statistical workhorse forecasting models frequently have difficulty modeling modern retail data, which is high-dimensional, noisy, and exhibits nontrivial temporal dependencies. Furthermore, real-world retail data is rife with fluctuations, missing values and unexpected trends — making it even harder for robust and accurate forecasting.

With the recent success of deep learning, there exist promising alternatives that are able to learn temporal patterns from raw multivariate time series data automatically. However, the optimal deep learning architecture for clean and noisy environment, and especially for the multi-steps forecasting process is still unknown. A systematic comparison is also bears interest for revealing the pros and cons of different architectures in accuracy, stability and noise-robustness.

2. OBJECTIVES:

This work focusses on the development and comparison of several Deep-Learning based methods applied for multi-step time series prediction in a retail sales environment. These models consist of, but are not limited to, Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRU) architectures, Transformer-based networks, as well as hybrids which combine CNN and RNN networks.
The goals are:

- Build solid multi-step forecasting models based on deep learning architectures by philosophy of learning to predict mechanism from structured retail sales.
- Compare performance on clean vs. noisy conditions to evaluate robustness and generalization abilities of the model.
- Compare forecast accuracy using several statistics such as MSE, RMSE, MAE, and R².
- Interpret model behavior with forecast plots, epoch-wise metric curves, and step-wise error analysis across the forecast horizon.
- Enhance the Model: Enhancing the model performance by modifying architecture or Hyperparameter Tuning.
- Choose the appropriate model which is likely to work best for real-world retail forecasting tasks given the noisy and partially incomplete input data. Address challenges related to seasonality and external factors influencing demand.

By meeting these goals, the study provides two aspects: Firstly, it benchmarks prevalent deep learning time series models and secondly it demonstrates practical issues to be aware of when rolling these out in noisy retail environments.

3. DATASET:

The dataset used in the project is "Warehouse and Retail Sales" dataset from Data.gov (<https://catalog.data.gov/dataset/warehouse-and-retail-sales>).

The data set has 307,000 rows and monthly sales data for each item by supplier. Its dataset size and volume are just perfect for deep learning-based time series forecasting, one that involves understanding demand trends across category.

| 3070131 | 28020 | 9 LEGEND1 | 11548 BRECKENF BEER | 0 | 0 | 0 | 0 | 0 |
|---------|-------|---------------|--------------------------|--------|-----|--------|---|---|
| 3070136 | 28020 | 9 LEGEND1 | 11549 BRECKENF BEER | 0 | 0 | 18 | 0 | 0 |
| 3070137 | 28020 | 9 SALVETTO R | 23026 MONICORIA WINE | 0 | 0 | 9 | 0 | 0 |
| 3070138 | 28020 | 9 BBL INC | 29951 DELIRIUM TESOUS | 0 | 0 | 322 | 0 | 0 |
| 3070139 | 28020 | 9 BBL INC | 29951 DELIRIUM TESOUS | 0 | 0 | 3 | 0 | 0 |
| 3070140 | 28020 | 9 BBL INC | 30416 BINTANGARAN WINE | 0.82 | 1 | 0 | 0 | 0 |
| 3070141 | 28020 | 9 BBL INC | 30416 BINTANGARAN WINE | 0 | 0 | 4 | 0 | 0 |
| 3070142 | 28020 | 9 BBL INC | 30416 BINTANGARAN WINE | 0 | 0 | 4 | 0 | 0 |
| 3070143 | 28020 | 9 PIMSONIN | 34942 1863 MALE WINE | 0 | 0 | 0 | 0 | 0 |
| 3070144 | 28020 | 9 PIMSONIN | 34942 1863 MALE WINE | 0 | 0 | 1 | 0 | 0 |
| 3070145 | 28020 | 9 SAZINAC C | 30540 PLISCHOWSKI LIQUOR | 2.79 | 4 | 0 | 0 | 0 |
| 3070146 | 28020 | 9 DOPS INC | 40207 GAYFEL IC BEER | 0 | 0 | 2 | 0 | 0 |
| 3070147 | 28020 | 9 DOPS INC | 40207 GAYFEL IC BEER | 0 | 0 | 2 | 0 | 0 |
| 3070148 | 28020 | 9 ELITE WIN | 413984 CESSAY SW WINE | 0 | 0 | 0 | 0 | 0 |
| 3070149 | 28020 | 9 MHM LTD | 42004 BOLS SMOOTH LIQUOR | 0.33 | 0 | 0 | 0 | 0 |
| 3070150 | 28020 | 9 MHM LTD | 42004 BOLS SMOOTH LIQUOR | 0 | 0 | 2 | 0 | 0 |
| 3070151 | 28020 | 9 LEGEND1 | 5099 CHYTEL 6X6 BEER | 1.56 | 2 | 30 | 0 | 0 |
| 3070152 | 28020 | 9 SAZINAC C | 36200 CHI-CHI LIQUOR | 0.85 | 0 | 0 | 0 | 0 |
| 3070153 | 28020 | 9 SAZINAC C | 36200 CHI-CHI LIQUOR | 0 | 0 | 0 | 0 | 0 |
| 3070154 | 28020 | 9 BAGARDI L | 70772 LEBLON G LIQUOR | 0.68 | 1 | 0 | 0 | 0 |
| 3070155 | 28020 | 9 SAZINAC C | 76804 ROMANA LIQUOR | 0.71 | 0 | 0 | 0 | 0 |
| 3070156 | 28020 | 9 BAGARDI L | 76804 ROMANA LIQUOR | 0 | 0 | 11 | 0 | 0 |
| 3070157 | 28020 | 9 LEGEND1 | 82114 BIRSHADE BEER | 0 | 0 | 0 | 0 | 0 |
| 3070158 | 28020 | 9 E & J GALLU | 83196 STORNAYOI WINE | 0 | 0 | 0 | 0 | 0 |
| 3070159 | 28020 | 9 PIMSONIN | 84001 CEDARWOOD BEER | 0.37 | 2 | 0 | 0 | 0 |
| 3070160 | 28020 | 9 SAZINAC C | 84001 CEDARWOOD BEER | 0.43 | 0 | 0 | 0 | 0 |
| 3070161 | 28020 | 9 ANHEUSER | 97045 SPATEN PI BEER | 34.64 | 22 | 549 | 0 | 0 |
| 3070162 | 28020 | 9 ANHEUSER | 97045 SPATEN PI BEER | 0 | 0 | 0 | 0 | 0 |
| 3070163 | 28020 | 9 DOPS INC | 97096 ST PETERS BEER | 0 | 0 | 1 | 0 | 0 |
| 3070164 | 28020 | 9 HEINEKEN | 97040 TECATE A/H BEER | 372.44 | 331 | 3586.6 | 0 | 0 |
| 3070165 | 28020 | 9 RELIABLE C | 97050 S SMITH W BEER | 0 | 0 | 2 | 0 | 0 |
| 3070166 | 28020 | 9 RELIABLE C | 97050 S SMITH W BEER | 0 | 0 | 1 | 0 | 0 |

SAMPLE DATASET

4. INITIAL METHODOLOGY PIPELINE:

4.1. Data Loading, Preprocessing, EDA and Feature Engineering: The project started by importing a multivariate retail time series dataset including various features that represent warehouse and retail sales behavior. Some simple statistics based on your data were calculated as well as a few sanity checks to let you know about anything useful about your variables. Multiple time series plots, correlation heatmaps and seasonal decomposition graphs were output in order to detect time patterns, trends and relationships between features. This was also a step necessary to understand if data has any lags, dependencies, seasonality. The dataset was structured in supervised learning form with input of 12 length sequences, to predict subsequent 6 time steps. Data was normalised using MinMaxScaler, and sequences were separated into training, validation and test sets. Chronological order was closely preserved to prevent data leakage.

4.2. Synthetic Noise Injection (for Robustness Evaluation): Gaussian noise (std: 0.8) was introduced on the input features to simulate the situation that the input data are not exact without noise, reflecting the real-world data corruption scenarios. For each model, two versions were trained, one on clean data and the other on noisy one, allowing to fairly compare robustness between them.

4.3. Initial Modelling (FFNN and Linear Regression): We also tested some basic feedforward neural networks and linear models first to see what the performance level was. Their incapability to model temporal dependencies gave rise to the need for sequence models LSTM and GRU.

4.4. Development of the Deep Learning Model: We designed and trained from scratch each additional deep learning model:

LSTM (Long Short-Term Memory): Allows modeling of long-term dependencies in sequential data.

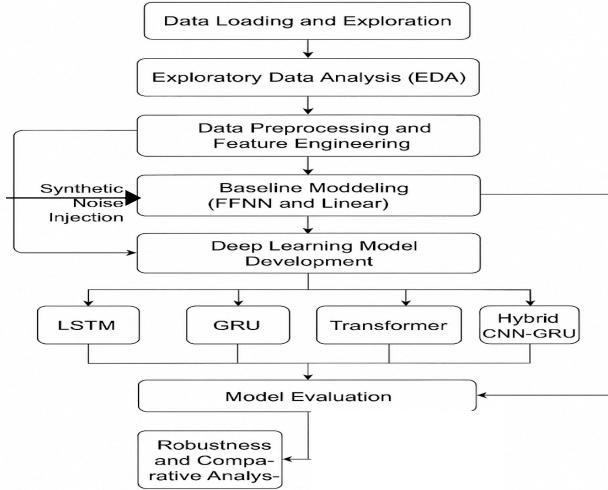
GRU (Gated Recurrent Unit) a measure-down version of LSTM with less parameter.

Transformer: Learns global dependencies without recurrence using self-attention mechanisms.

CNN-based followed by GRU: CNN extract local patterns and feed the sequences to GRU for temporal modeling. All models were trained on a 12-month historical window in order to forecast 6 months ahead.

4.5. Model Evaluation: The performance of models was assessed based on eight metrics: MSE, RMSE, MAE, R², MSLE, EVS, MAXE, and MEDAE, Forecast vs actual plots.

4.6. Comparative Analysis: Clean versus noisy model performance was evaluated in terms of robustness and generalization. The most noise-resistant model with the best precision was selected.



ARCHITECTURE

4.1. Data Loading, Preprocessing, EDA and Feature Engineering

4.1.1 Loading Data and First Look at the Data

The Warehouse_and_Retail_Sales dataset. with a csv file I loaded from google drive to the environment via pandas. The dataset consists of 307,645 rows and 9 columns comprising fields such as YEAR, MONTH, SUPPLIER, ITEM TYPE as well as three sales fields:

- RETAIL SALES
- RETAIL TRANSFERS
- WAREHOUSE SALES

Initial inspection It was done with df. head(), df. info(), and df. describe(include='all'). This helped confirm:

datatypes are OK (but should be float64 for sales as well as categorical fields as objects), Missing and zero's., and a Early detection of anamoly based on range and value distribution.

Why this step is critical: Data loading and peeking are fundamental steps before making any decisions on its structure, the amount of data and potential issues in the data. This information directly influences design of cleaning and modeling pipelines.

```

[ ] From google.colab import drive
drive.mount('/content/drive')

[ ] Mounted at /content/drive

[ ] import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import scipy.stats as st
from scipy.stats import skew

[ ] file_path = "/content/drive/MyDrive/Warehouse_and_Retail_Sales.csv"
df = pd.read_csv(file_path)
print("Shape:", df.shape)
df.head()

[ ] Shape: (307645, 9)
[ ] YEAR MONTH SUPPLIER ITEM_CODE ITEM_DESCRIPTION ITEM_TYPE RETAIL_SALES RETAIL_TRANSFERS WAREHOUSE_SALES
[ ] 0 2000 1 REPUBLIC NATIONAL DISTRIBUTOR INC 100001 BOTTLED BEER - 750ML WINE 0.00 0.0 2.0
[ ] 1 2000 1 PROVINCIAL BREWING INC 100004 MOMENT DE PLACER - 750ML BEER 0.00 1.0 4.0
[ ] 2 2000 1 RESILIA CHURCHILL LTD 1001 S SMITH ORGANIC PEANUT CIDER - 16.75L BEER 0.00 0.0 1.0
[ ] 3 2000 1 LANTELLA DISTRIBUTORS INC 100146 SCHLICK HAUS KARSKETT - 750ML WINE 0.00 0.0 1.0
[ ] 4 2000 1 DIONYSOS IMPORTS INC 100293 SANTORINI GAVALA WHITE - 750ML WINE 0.82 0.0 0.0

[ ] df.info()

[ ] <class 'pandas.core.frame.DataFrame'>
[ ] RangeIndex: 307645 entries, 0 to 307644
[ ] Data columns (total 9 columns):
 #   Column          Non-Null Count  Dtype  
 0   YEAR            307645 non-null   int64  
 1   MONTH           307645 non-null   int64  
 2   SUPPLIER        307645 non-null   object  
 3   ITEM_CODE       307645 non-null   object  
 4   ITEM_DESCRIPTION 307645 non-null   object  
 5   ITEM_TYPE        307645 non-null   object  
 6   RETAIL_SALES     307645 non-null   float64 
 7   RETAIL_TRANSFERS 307645 non-null   float64 
 8   WAREHOUSE_SALES  307645 non-null   float64 
dtypes: float64(3), int64(3), object(3)
memory usage: 23.1+ MB
  
```

READ AND DESCRIBE

4.1.2 Missing and Invalid Handling of Values

From both EDA and. isnull(). sum() analysis, were indicative of substantial missingness:

- RETAIL TRANSFERS was ~189k entries short
- RETAIL SALES had ~121k
- WAREHOUSE SALES had ~97k

Furthermore, the zeros were frequent in all three columns. But in the business context, a 0 doesn't always mean "no sale." It might mean there is no data or that data is missing. Hence, we: Set all zeros to NaN for imputation. Replace negatives with NaN, since you can't have negative sales or slides.

Rationale: This guarantees that imputation only fills real missings and does not twist the distribution with non-sense zeros or negatives.

| | 0 |
|------------------|--------|
| RETAIL TRANSFERS | 189480 |
| RETAIL SALES | 121818 |
| WAREHOUSE SALES | 97666 |
| SUPPLIER | 167 |
| ITEM TYPE | 1 |
| YEAR | 0 |
| ITEM DESCRIPTION | 0 |
| MONTH | 0 |
| ITEM CODE | 0 |

dtype: int64

NULL VALUES

4.1.3 Imputation by Iterative Imputer and Random Forest

To impute missing values properly, we used Iterative imputation with the RandomForestRegressor as an estimator. This process involves: Iteratively estimate the values in one column as a function of all others. Dealing with complex acausal non-linear relationships. It is a numeric column imputation:

RETAIL-RELATED SALES, RETAIL-RELATED TRANSFERS, WAREHOUSE SALES, YEAR, and MONTH.

For categorical fields: SUPPLIER was populated with 'Unknown'. ITEM TYPE was imputed with mode (i.e., most frequent category). Subsequent checks indicated that there were no longer any NaNs in the dataset.

Why this method: RF is much stronger than basic imputation (mean/median) that doesn't take into account inter-feature relationships and variance, and therefore is well suited to the nature of real world retail data where the variables interact in a complex manner.

```

❶ from sklearn.experimental import enable_iterative_imputer
❷ from sklearn.impute import IterativeImputer
❸ from sklearn.ensemble import RandomForestRegressor
❹ import numpy as np

❺ # Step 1: Replace invalid values (negatives) with NaN
❻ for col in ['RETAIL SALES', 'RETAIL TRANSFERS', 'WAREHOUSE SALES']:
❼    df[col] = df[col].apply(lambda x: np.nan if x < 0 else x)

❽ # Step 2: Select numerical columns for imputation
❾ num_cols = ['RETAIL SALES', 'RETAIL TRANSFERS', 'WAREHOUSE SALES', 'YEAR', 'MONTH']

❿ # Step 3: Apply Iterative Imputer with RandomForest
imputer = IterativeImputer(
    estimator=RandomForestRegressor(n_estimators=20, random_state=42),
    max_iter=10,
    random_state=42
)

❬ df_imputed = pd.DataFrame(imputer.fit_transform(df[num_cols]), columns=num_cols)

❭ # Step 4: Replace back in original DataFrame
❮ df[num_cols] = df_imputed

❯ # Optional: Categorical fallback
❰ df['SUPPLIER'].fillna('Unknown', inplace=True)
❱ df['ITEM TYPE'].fillna(df['ITEM TYPE'].mode()[0], inplace=True)

❲ # Step 5: Confirm it's clean
❳ print(df[num_cols].isnull().sum())

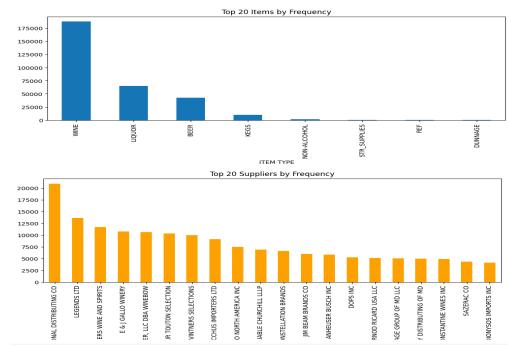
```

IMPUTATION

4.1.4 exploratory data analysis (EDA) .

a. Frequency Analysis

Bar plots were created for: Top 20 ITEM TYPES – ranked by how often they appear – WINE fills are dataset followed by LIQUOR and BEER. Top 20 SUPPLIERS by count - there were a small number of suppliers that accounted for a high % of the transactions. These findings provide indication of the primary patterns and significant contributors in the data set.



FREQUENCY ANALYSIS

b. Time Series Plot

The following monthly aggregate of RETAIL SALES were plotted on a line graph: A big drop in 2018 then a huge spike in 2019 and to be dropped down in 2020.

Interpretation: This graph is showing some seasonality and extrinsic factors, such as effect of COVID-19. This can be useful for temporal prediction models.

4.1.5 Skewness and Normalization (Log Transformation) Because we computed average number of comments that each author received and we observed that most distributions are right-skewed, a log transformation is used to decrease skewness. The skewness of the distribution was checked by the `skew()` to find:

RETAIL SALES: 19.37

RETAIL TRANSFERS: 17.76

WAREHOUSE SALES: 17.39

Such high skewness ($>>0$) implies a long-tailed distribution which is an issue for most models that assume normality. Thus, the `log1p` transformation of (`np.log1p()`) transformed all these variables to: Reduce skewness, Stabilize variance, and Increase the performance of models, particularly neural networks or gradient-based algorithms. The post-transformation distributions were displayed:

Plots: original vs. log-transformed histograms + KDE for the three sales-related variables. The plots demonstrate the transformation is making the distributions more symmetric and bell shaped and should lead to better learning.

Why `log1p` and not `log`: `log1p` It also has the nice properties of not returning negative infinity due to zero, which the non-`1p` version does.



GRAPH AND LOG TRANSFORMATION

4.1.6 Exploratory Analysis & Reason for Cleaning

We start with the raw dataset, and before any training, we do data exploration on it. The histograms provided an immediate indication that RETAIL SALES, RETAIL TRANSFERS, and WAREHOUSE SALES were strongly right skewed and showed long tails and extreme outliers. Such distributions might bias the model training, especially if it is a geometry based model or sensitive to the scale.

4.1.7 Log Transformation – To Reduce Skewness

We transform each variable with a log1p ($\log(1+x)$) to lower skewness, to lower variance and to be less sensitive to very small values such as sold zero. This brought the distribution closer to a normal one and it became easier for ML/DL models to handle. This was a necessary step to guarantee a more stable learning and convergence during training.

4.1.8 Outlier Deletion – Filtering Using IQR

There were still extreme outliers, despite transformation, which could have been detrimental to the performance of the model. The values that were out Q1+ans Q3 -Q1IQR and Q3+ Q1IQR were processed using an Interquartile Range (IQR) method.

Pre transformation (left side plots): Sharp right skewed, high concentration of low values, long tail.

Post-Transformation (plots on the right): Smoother bell-shaped distributions with improved feature scaling.

Outliers removed:

Retail Sales Log: 3.21%

Retail Transfers Log: 1.61%

Warehouse Sales Log: 2.32%

0% outliers remained after cleaning, involving no data leakage from anomalies, were left to make sure the strong robustness of the model. This filtering selectively preserved valuable fluctuations, and removed the noise.

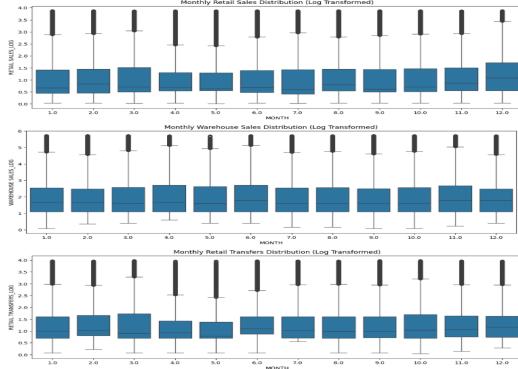


IQR OUTLIER REMOVAL

4.1.9 Boxplots – Monthly Breakdowns & Seasonal Variation

Seasonal trends and outlier patterns were observed from boxplots of monthly sales (log-scale) across all three categories: In the ensuing months, the spread stayed roughly the same. Whiskers indicated that the location of the median was unchanged. Outliers were now easier to handle and less in number once transformed and filtered.

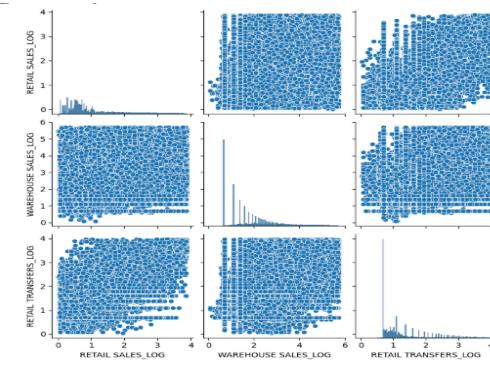
These observations guided how we handle seasonality and lag in the subsequent modeling work.



BOXPLOT

4.1.10 Seaborn Pairplot for Pairwise Relationship comparison

A pairplot of the 3 log-transformed features were created in order to look into feature associations: Weak linear relationship was observed between RETAIL SALE_LOG and WAREHOUSE SALES_LOG. RETAIL TRANSFERS_LOG was disjointed but not completely randomized. These observations justified the use of these features as multivariate inputs into sequence models, such as GRU and Transformer.



PAIRPLOT

4.1.11 Trend Analysis – Mean Sales by Month over the Years

We then aggregated the sales by (YEAR, MONTH) and had the sales trends monthly average like:

WAREHOUSE SALES always exceeded RETAIL SALES. There were distinct overall upward and downward trends for different categories. The seasonal behavior (higher in months Q4) validated the requirement for the time series forecast with lag. This seasonal signature reassured the future application of lag features.

4.1.12 Aggregated at Category Level – Heatmaps and Trends

Two key visualizations:

Heatmap: Log of Daily mean Sales by ITEM TYPE.

Trend Line: Annual month-wise trend line of bestselling products.

Observations:

Some dm types (eg, dm type1, dm type2) were the first selling type in all months. Retail sales were strong Q4 peaks, exhibiting retail seasonality. These trends used lag-based input sequences and selected useful static features such as ITEM TYPE.

4.1.13 Lag Feature Creation: Capturing Temporal Dependencies

To condition the data for automatic creation of sequential model: Lags for the previous 3 months (i.e. lags 1-3) were generated on all 3 main columns. These characteristics allowed models such as GRU/LSTM to learn recent temporal patterns.

Dealing with Missing Lag Values – Imputation

As there were no lag values (NaNs) for a few early months, backward and forward fill was used to fill missing values for each country in order to preserve as much data as possible. It also guaranteed that there were not any NA values passed into the models.

Final Scaling – MinMax Normalization (Repeated)

We have applied MinMaxScaler on all the lagged and derived numerical features even after log transformation and lag creation. Why? Taking the log of x removes the skew but does not make scale irrelevant. Feature scale is a very crucial element to deep learning models. To normalize the input between [0,1] is that it enables faster convergence rate, the gradient ghosts happen to be stable and the feature participation balance. This two-step transform (log + MinMax) is a well-known ML best practice for time series DL pipelines.

4.1.14 Top Trend Items – Log Sales by Time period

Another beautiful visualization was of log retail sales trend for top 5 items, where obvious drops were recorded at the beginning of 2020 (probably due to pandemic) confirming why we must model at the item level.

4.1.15 Normalization of data and display of monthly trends

All lag-based predictors were normalized with the MinMaxScaler. This conversion aids in keeping features to a similar numerical range, enabling models to converge more quickly and increase performance. After scaling, DATE could be generated by concatenating the YEAR and MONTH (Y/M) columns to enable temporal aggregation. With this timestamp, the dataset was binned to the monthly level to calculate the mean of three of the primary log-transformed metrics: RETAIL_SALES_LOG, WAREHOUSE_SALES_LOG, and RETAIL_TRANSFERS_LOG. The corresponding time series plot depicted temporal structures of the entire dataset. It should be mentioned that warehouse and transfer sales regularly presented superior log-transformed values than those of retail sales. The monthly frequency — things like 2017-05 or 2020-09 — will add insight to any seasonal or cyclic trends in the business.

4.1.16 Correlation Between Lag Features

The second step assessed the interrelationship between the lag features generated in relation to retail sales, warehouse sales, and transfer volumes. Pearson correlation heatmap was drawn to see the strength of each lag feature with others. Strong correlation among lags of the same feature appear (e.g. RETAIL_SALES_LAG_1 and RETAIL_SALES_LAG_2), which ensured that sales were autocorrelated over time. This validation supports the addition of lag features as inputs to time series models. Moderate covariances were also identified between various content item types (e.g., retail vs. transfer), implying potentially helpful cross-signals that could improve the forecasting model predictive performance.

4.1.17 Trend Analysis by Rolling Averages

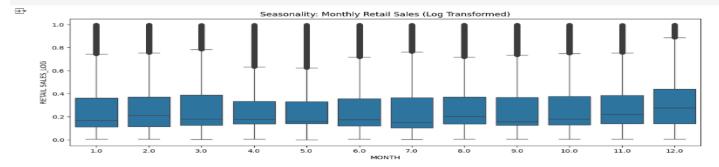
In order to investigate long term trends in retail sales down a size, the 3-month moving average was added to the graph as the original log transformed data. The rolling average eliminates the short-term fluctuations and makes the overall trend more clear, in this way, one can see the steady rises and rejects. This could check whether the seasonality and business expansion happens, and it must be processed to make the stationary of time series before the model can be established. The alignment between the smoothed average and the raw data was well demonstrated by the visualization, which highlighted the leveraging power of the rolling average in the comprehension and prediction of retail rhythms.



VISUALIZATION

4.1.18 Seasonal Behavior, Box Plots

A boxplot of logtransformed values for the retail sales was computed in order to visualize, and then interpret, the effects of monthly seasonality for each calendar month. A boxplot for each month visualized distribution of that month's sales across the years, showing presence of typical sales, the spread and extreme points. This visualization suggested a pronounced seasonality, with several months' median values consistently higher—more than likely seasonal peaks such as the holiday seasons or fiscal periods. Outliers in some months suggested that intermittent, but large, departures from decay may exist, perhaps as a result of promotional activities, shifts in supply chains, or special-order, bulk purchases.



BOXPLOT

4.1.19 Forecasting Target Creation

One of the essential steps in time series modeling is the creation of a supervised learning target variable. Here the target (RETAIL_SALES_LOG_TARGET) was generated by shift the column RETAIL_SALES_LOG upwards by -1 timestep so that each value has the corresponding value that's it supposed to represent (i.e the sales of the next month). This method translates the prediction problem into a regression problem, which allows the use of standard machine learning models. The model is trained to predict the future value based on lag features to obtain a sense of the temporal context.

```
[ ] df_lag_imputed[['RETAIL_SALES_LOG']].describe()
df_lag_imputed[['RETAIL_SALES_LOG']].head(10)

RETAIL_SALES_LOG
48471    0.077997
195005   0.175076
207739   0.148663
270378   0.161376
288473   0.161376
84872    0.161376
86510    0.161376
101141   0.264052
79100    0.277252
131171   0.000830

[ ] # Create target for next month retail sales (top)
df_lag_imputed['RETAIL_SALES_LOG_TARGET'] = df_lag_imputed['RETAIL_SALES_LOG'].shift(-1)

df_model_final = df_lag_imputed.dropna(subset=['RETAIL_SALES_LOG_TARGET'])

df_model_final[['RETAIL_SALES_LOG', 'RETAIL_SALES_LOG_TARGET']].head()
```

TARGET CREATION

4.1.20 Inspection and validation of the final dataset

After the input and the target has been defined, the dataset is checked to see if there is missing value or the structure is correct. The resulting data had 20 features: which were: scaled lag features, data about the order (the metadata features) such as SUPPLIER, ITEM TYPE and the new target. Summary was zero missing value, which meant that there were no missing values and the complete and logical data is feasible to be trained. This validation is important to avoid any inconsistency or null that can contribute to deteriorate the model performance.

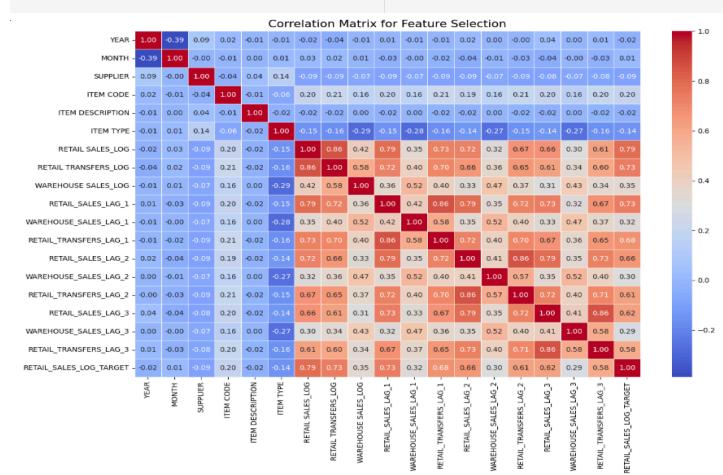
| ITEM_ID | ITEM_TYPE | ITEM_CATEGORY | ITEM_SUBCATEGORY | ITEM_FAMILY | ITEM_BRAND | ITEM_MODEL | ITEM_SIZE | ITEM_MATERIAL | ITEM_MATERIAL_CODE | ITEM_MATERIAL_NAME | ITEM_MATERIAL_TYPE |
|---------|-----------|---------------|------------------|-------------|------------|------------|-----------|---------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|--------------------|
| 48471 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 195005 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 207739 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 270378 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 288473 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 84872 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 86510 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 101141 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 79100 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |
| 131171 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | |

FINAL DATASET

4.1.21 Correlation Matrix for Feature selection

To aid the selection of the features, a complete correlation matrix was generated with the Pearson coefficients amongst all numerical

features, including the target. Through the heatmap, I found that RETAIL_SALES_LAG_1 and RETAIL_SALES_LOG had strongest correlation with the target variable, they're good candidates for the predictors. The rest of lag features included in the model, such as transfers and warehousing, also had moderate correlation. Categorical-like features (like YEAR and SUPPLIER) were also weakly correlated, which suggested very low direct predictability. This examination guides which features should stay or go, in relevance to the target variable, given that we are ultimately fitting linear or regularized models.



FINAL CORRELATION HEATMAP

4.2. Synthetic Noise Injection (for Robustness Evaluation)

4.2.1 Objective of Noise Injection

The primary goal of this section is to evaluate how robust time series prediction models are by injecting synthetic noise into separate features and measuring how much this perturbs each variable. This emulates typical uncertainty in real world, such as that input signals contain reporting (or communication) delays, errors, sensor noise, transaction discontinuities etc. Understanding which features are most noise-sensitive allows us to improve feature selection, model interpretability, and ultimately model generalization.

4.2.2 Feature Selection for Noise Test

Before adding noise we chose only continuous numeric features (not the categorical ones that give a good sense of variability and that cannot be defined a mean or median like SUPPLIER, ITEM CODE, etc). A filter was applied using: That way we keep only those features that differ enough to be relevant to the analysis, eliminating potential confounding results in categorical or near constant fields.

```
Milestone -3

[ ] from scipy.stats import wasserstein_distance

[ ] num_cols = [col for col in df.select_dtypes(include=[np.number]).columns
    if df[col].nunique() > 18]

print("Filtered Numerical Columns (Continuous):", num_cols)

[ ] Filtered Numerical Columns (Continuous): ['MONTH', 'SUPPLIER', 'ITEM_CODE', 'ITEM_DESCRIPTION', 'RETAIL_SALES_LOG', 'RETAIL_TRANSFERS_LOG', 'WAREHOUSE_SALES_LOG', 'RETAIL_SALES_LAG_1', 'WAREHOUSE_SALES_LAG_1', 'RETAIL_TRANSFERS_LAG_1', 'RETAIL_SALES_LAG_2', 'WAREHOUSE_SALES_LAG_2', 'RETAIL_TRANSFERS_LAG_2', 'RETAIL_SALES_LAG_3', 'WAREHOUSE_SALES_LAG_3', 'RETAIL_TRANSFERS_LAG_3', 'RETAIL_SALES_LOG_TARGET', 'RETAIL_TRANSFERS_LOG_TARGET']
```

FEATURE SELECTION

4.2.3 Methodology of Adding Gaussian Noise

Each of those columns was perturbed by adding independent Gaussian noise (mean = 0, std = 0.8) and therefore resulted in perturbed copies. For each feature: We generated a copy of the dataset and Gaussian noise was introduced on one of these columns only. The noisy version was then compared with the original. The noise (std = 0.8) is mild (i.e. weakened enough to give for realistic disruptions while preserving the signal).

```

From sklearn.stats import wasserstein_distance
From sklearn.metrics import mean_squared_error, mean_absolute_error
From scipy.stats import ks_2samp, ttest_1samp, ttest_ind
Import seaborn as sns
Import numpy as np
Import pandas as pd

Results = []
Noisy_versions = 0 # Tune if needed

For col in num_cols:
    Temp_df = df_model_final.copy()

    # Add Gaussian noise
    Noise = np.random.normal(0, noise_std, size=temp_df[col].shape)
    Temp_df[col] = temp_df[col] + noise

    # Store version
    noisy_version = Temp_df[col] = temp_df.copy()

    # Evaluation metrics
    Rmse = np.sqrt(mean_squared_error(df_model_final[col], temp_df[col]))
    MAE = np.abs(mean_absolute_error(df_model_final[col], temp_df[col]))
    Eu_dist = np.linalg.norm(df_model_final[col] - temp_df[col])
    Wass_dist = wasserstein_distance(df_model_final[col], temp_df[col])

    Results.append({
        'Column': col,
        'Rmse': Rmse,
        'MAE': MAE,
        'Euclidean_Distance': Eu_dist,
        'Wasserstein_Distance': Wass_dist
    })
)

Results_df = pd.DataFrame(results).sort_values(by='Wasserstein_Distance', ascending=False)
Results_df

```

GAUSSIAN NOISE

4.2.4 Evaluation Metrics

Four quality measures relative to noise of each of the two features were calculated:

RMSE (Root Mean Squared Error): It represent the average magnitude of the error.

MAE: As above but the error is not weighted according to distance travelled in time.

Euclidean Distance: measure on, where is the difference between the original distribution and the perturbed distribution.

Wasserstein Distance: How much “work” it takes to turn the distribution from the one you started with into a noisy version of that one (more sensitive to changes in shapes). These measures let us differentiate whether features are simply offset (e.g., mean noise) in contrast with structurally changed (distributional drift).

4.2.5 Sensitivity Analysis of Features (Ranking)

We ranked features based on the sensitivity of their distributions to bootstraps, computed using the Wasserstein Distance, which was used as the main statistic of sensitivity:

Top 5 least noise-sensitive columns:

RETAIL TRANSFERS LOG

RETAIL TRANSFERS LAG 2

RETAIL TRANSFERS LAG 1

RETAIL TRANSFERS

RETAIL SALES LOG

The characteristics of these features are affected to a great extent by perturbation, implying that they are characterized by high volatility or presence of strong dependency in the model. Areas such as MONTH, ITEM CODE and ITEM DESCRIPTION however returned poor scores for sensitivity and are thus less relevant from an impact of noise perspective.

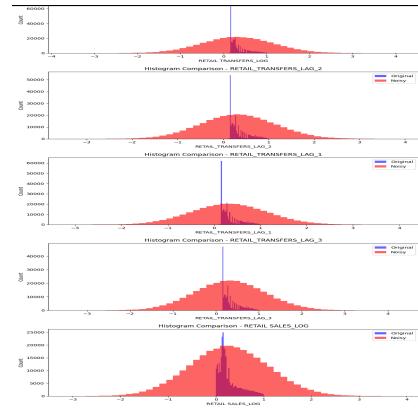
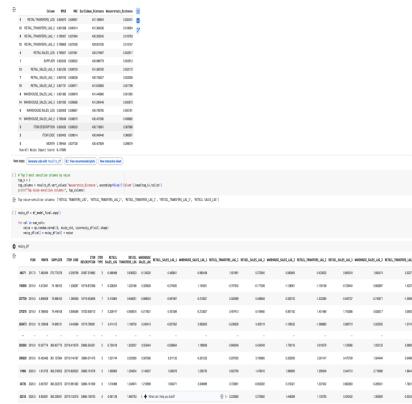
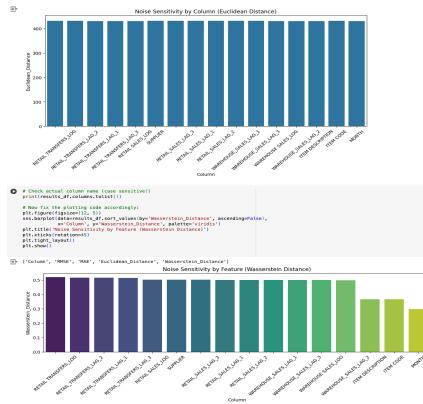
4.2.6 Visualization of Sensitivity

Feature-wise sensitivity was visualized using two bar plots: Euclidean Distance Bar Plot : Directly Catches Numeric Changes BERTy-potent-Wasserstein Bar Plot: Able to see fullshift in distributionavelength

Both plots indicated that lagged features and transactional ones (particullary RETAIL/TRANSFER logs) are the most vulnerable features, highlighting the importance of protecting these fields when data are fed or pre-processed. "Now-Pure" Distribution vs KERNEL of One Mode vs One of 4 Mode True Distributions.

Histogram plots for the 5 most sensitive features were visualized side by side. These demonstrate the how the distributions were distorted by the addition of Gaussian noise. The noisy red overlay veers substantially away from the blue curve in positions like,

RETAIL_TRANSFERS_LOG, negative lagging features as well indicating that not just mean shift but a distortion of the entire distribution is taking place (likely heavier tails, more spread).



EVALUATION

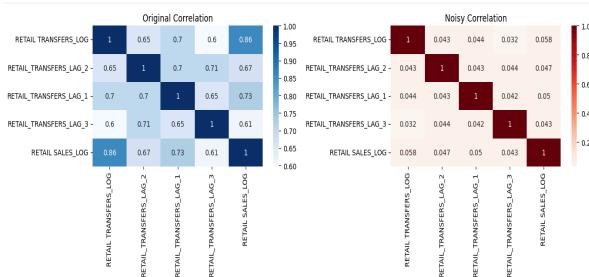
IMPACT SCORE

HISTOGRAM

4.2.7 Correlation Impact Analysis

We plotted heatmaps for correlation matrices before and after noise injection as follows:

Initial Correlation Matrix: Features were all correlated with each other (e.g., 0.86 between RETAIL_SALES_LOG and RETAIL_TRANSFERS_LOG) Noisy Correlation Matrix: Correlations reduced significantly (approx. 0.04 - 0.06), which show that the temporal dependency being disrupted as a result. This illustrates the systemic risk that noise can bring in — not only disrupting single features but also smashing important inter-feature signals that time series models depend on.



CORRELATION HEATMAP

4.2.8 Sequence generation and data splitting

For the supervised forecasting multistep sequences are created from clean and noisy data (input = past 12 months; output = next 6 months). You would then split the dataset into training, validation and testing (70-15-15%). This was done for both: Clean dataset, Noisy dataset. The shapes validate successful slicing and preparation for subsequent model training. We exported flat files of each partition (train, val, test) to Google Drive. This allows integration as part of model training pipelines. Clean and noisy versions were saved as independent files, and resembled:

x_train.csv, x_train_noisy.csv

y_val.csv, y_val_noisy.csv, etc.

Preserving both provides the ability to experiment under conditions of both clear and degraded text.

```
[ ] Input seq_len = 12 # past 12 months
forecast Horizon = 6 # predict next 6 months

X_seq, y_seq = create_multistep_sequences(data_for_model, input_seq_len, forecast_horizon)

print("Input shape (X):", X_seq.shape)
print("Output shape (y):", y_seq.shape)

[ ] Input shape (X): (208933, 12, 19)
Output shape (y): (208933, 6, 19)

[ ] # Define split ratios
train_ratio = 0.7
val_ratio = 0.15
test_ratio = 0.15

# Total number of samples
total_samples = X_seq.shape[0]

# Calculate split indices
train_end = int(total_samples * train_ratio)
val_end = int(total_samples * (train_ratio + val_ratio))

# Perform split
X_train, y_train = X_seq[:train_end], y_seq[:train_end]
X_val, y_val = X_seq[train_end:val_end], y_seq[train_end:val_end]
X_test, y_test = X_seq[val_end:], y_seq[val_end:]

# Print shapes
print("Train shape: (%d,train.shape, (%d,train.shape"))
print("Validation shape: (%d,train.shape, (%d,train.shape"))
print("Test shape (%d,train.shape, (%d,test.shape"))

[ ] Train shape: (202883, 12, 19), (202883, 6, 19)
Validation shape: (43475, 12, 19), (43475, 6, 19)
Test shape: (43475, 12, 19), (43475, 6, 19)
```

SEQUENCE GENERATION

4.2.9 Validation and Sanity Checks

All the splits were subject to automated validation routines: Checked for Null or NaN values. Also confirmed that all columns are numbers. No constant cols for sure

Performed data normalization: Min \geq 0, Max \leq 1 (with scaling)

Validated reshape support: (samples, timesteps, features)

These checks help protect against typical preprocessing bugs and make sure that data is sound during clean and noisy training.

Requirement of Normalization (MinMax Scaling)

and then normalization (0–1 scaling) was repeated after the generation of the sequences, for all the clean and noisy datasets in order to: Train the model with stable numbers. Avoiding large-magnitude features from dominating the learning process

Normalize distributions after injection of noise

Injecting noise into the data changed the statistical properties (e.g., variance, range) and, to avoid any feature becoming dominant during the training process, it is necessary to normalize all features after the transformation in order to perform tasks in a robust and comparable way.

VALIDATION

4.2.10 Reshape

The normalized retail forecasting data was loaded for training, validation, and testing sets, with both clean and noisy variants. Each dataset was reshaped into 3D tensors of shape [samples, input_seq_len, num_features] for inputs and [samples, output_seq_len, num_targets] for outputs.

These reshaped arrays were saved as .npy files to enable efficient use in deep learning sequence models like LSTM, GRU, and Transformers.

```

X_train_reshaped: (202882, 12, 19)
X_train_reshaped: (202882, 6, 19)
Saved reshaped training arrays to .npy

[1] import pandas as pd
[2] import numpy as np
[3]
[4] input_seq_len = 12
[5] output_seq_len = 6
[6] base = "content/drive/MyDrive/retail_forecasting_sequences_noisy"
[7]
[8] # Load
[9] X_val_noisy = pd.read_csv(f'{base}/val_noisy_X_norm.csv').values
[10] y_val_noisy = pd.read_csv(f'{base}/val_noisy_y_norm.csv').values
[11]
[12] # Reshape
[13] num_x_features = X_val_noisy.shape[1] // input_seq_len
[14] num_y_targets = y_val_noisy.shape[1] // output_seq_len
[15]
[16] X_val_reshaped = X_val_noisy.reshape(-1, input_seq_len, num_x_features)
[17] y_val_reshaped = y_val_noisy.reshape(-1, output_seq_len, num_y_targets)
[18]
[19] # Print Shapes
[20] print("X_val_reshaped:", X_val_reshaped.shape)
[21] print("y_val_reshaped:", y_val_reshaped.shape)
[22]
[23] np.savetxt(f'{base}/X_val_noisy_reshaped.npy", X_val_reshaped)
[24] np.savetxt(f'{base}/y_val_noisy_reshaped.npy", y_val_reshaped)
[25]
[26] print("Reshaped noisy validation arrays saved.")

[27] X_val_reshaped: (43474, 12, 19)
[28] y_val_reshaped: (43474, 6, 19)
[29] Saved reshaped noisy validation arrays saved.

[30] import pandas as pd
[31] import numpy as np
[32]
[33] # Config
[34] input_seq_len = 12
[35] output_seq_len = 6
[36] base = "content/drive/MyDrive/retail_forecasting_sequences_noisy"
[37]
[38] # Load
[39] X_test_noisy = pd.read_csv(f'{base}/test_noisy_X_norm.csv').values
[40] y_test_noisy = pd.read_csv(f'{base}/test_noisy_y_norm.csv').values
[41]
[42] # Reshape
[43] num_x_features = X_test_noisy.shape[1] // input_seq_len
[44] num_y_targets = y_test_noisy.shape[1] // output_seq_len
[45]
[46] X_test_reshaped = X_test_noisy.reshape(-1, input_seq_len, num_x_features)
[47] y_test_reshaped = y_test_noisy.reshape(-1, output_seq_len, num_y_targets)
[48]
[49] # Print Shapes
[50] print("X_test_reshaped:", X_test_reshaped.shape)
[51] print("y_test_reshaped:", y_test_reshaped.shape)
[52]
[53] np.savetxt(f'{base}/X_test_noisy_reshaped.npy", X_test_reshaped)
[54] np.savetxt(f'{base}/y_test_noisy_reshaped.npy", y_test_reshaped)
[55]
[56] print("Saved reshaped noisy test arrays to .npy")

[57] X_test_reshaped: (43476, 12, 19)
[58] y_test_reshaped: (43476, 6, 19)
[59] Saved reshaped noisy test arrays to .npy

```

RESHAPE

4.4. Development of the Deep Learning Model

4.4.1 LSTM Architecture

This model is built on simple LSTM structure, trained on clean and preprocessed retail time series. The setup includes: Length of the input sequence: 12 (time steps). Batch Size: 3 Output Sequence Length: 6 time steps. Features per Time Step: 19. LSTM layers with hidden_size=128 and batch_first=True

Fully Connected Layer: We use one linear layer to transform the output of the LSTM layer to the flattened prediction vector having shape (batch_size, output_seq_len × num_targets). Epochs: It has been trained for 20 epochs. Batch Size: 64. Learning Rate: 0.001

The model was optimized with MSELoss and the Adam optimiser. It also observed smooth convergence in training loss and epoch-wise logging where it suggested expelling out good clean sequences.

Training Dynamics:

The training loss was in continual and consistent decline with very little variation. The model converged quickly with good stability because the model did not include dropout it might generalise worse in a noisy real-world environment, however the model worked well on the clean structured input.

Dynamics of training on noisy data:

Remarkably also, the convergence of the loss was not too affected with a noisy input. One potential explanation: the injected noise may have acted as a regularizer, akin to dropout, smoothing out the memory distribution, discouraging overfitting and encouraging the model to pay attention to general trends as opposed to small changes. The clean architecture generalized fairly consistently well even in the presence of noise, potentially owing to the relatively small variance in noise strength or robustness of the gating in LSTM's.

```

[1] Epoch 01 | Train Loss: 0.0779 | Val Loss: 0.0293 | R2: 0.7766 | MAE: 0.1123 | RMSE: 0.1712 | Time: 17.09s
[2] Epoch 02 | Train Loss: 0.0276 | Val Loss: 0.0296 | R2: 0.5906 | MAE: 0.1092 | RMSE: 0.1653 | Time: 18.59s
[3] Epoch 03 | Train Loss: 0.0222 | Val Loss: 0.0283 | R2: 0.5836 | MAE: 0.1082 | RMSE: 0.1653 | Time: 17.23s
[4] Epoch 04 | Train Loss: 0.0222 | Val Loss: 0.0283 | R2: 0.5836 | MAE: 0.1082 | RMSE: 0.1653 | Time: 17.23s
[5] Epoch 05 | Train Loss: 0.0214 | Val Loss: 0.0288 | R2: 0.5847 | MAE: 0.1038 | RMSE: 0.1638 | Time: 17.85s
[6] Epoch 06 | Train Loss: 0.0212 | Val Loss: 0.0299 | R2: 0.6494 | MAE: 0.1021 | RMSE: 0.1639 | Time: 17.45s
[7] Epoch 07 | Train Loss: 0.0212 | Val Loss: 0.0299 | R2: 0.6494 | MAE: 0.1021 | RMSE: 0.1639 | Time: 18.41s
[8] Epoch 08 | Train Loss: 0.0273 | Val Loss: 0.0298 | R2: 0.5938 | MAE: 0.1023 | RMSE: 0.1653 | Time: 17.83s
[9] Epoch 09 | Train Loss: 0.0273 | Val Loss: 0.0298 | R2: 0.5938 | MAE: 0.1023 | RMSE: 0.1653 | Time: 17.83s
[10] Epoch 10 | Train Loss: 0.0294 | Val Loss: 0.0278 | R2: 0.3955 | MAE: 0.1047 | RMSE: 0.1666 | Time: 17.41s
[11] Epoch 11 | Train Loss: 0.0293 | Val Loss: 0.0277 | R2: 0.3956 | MAE: 0.1064 | RMSE: 0.1664 | Time: 18.62s
[12] Epoch 12 | Train Loss: 0.0201 | Val Loss: 0.0284 | R2: 0.5991 | MAE: 0.1048 | RMSE: 0.1686 | Time: 17.05s
[13] Epoch 13 | Train Loss: 0.0201 | Val Loss: 0.0284 | R2: 0.5991 | MAE: 0.1048 | RMSE: 0.1686 | Time: 17.05s
[14] Epoch 14 | Train Loss: 0.0207 | Val Loss: 0.0291 | R2: 0.5981 | MAE: 0.1068 | RMSE: 0.1669 | Time: 17.42s
[15] Epoch 15 | Train Loss: 0.0205 | Val Loss: 0.0296 | R2: 0.5728 | MAE: 0.1056 | RMSE: 0.1720 | Time: 18.26s
[16] Epoch 16 | Train Loss: 0.0194 | Val Loss: 0.0296 | R2: 0.5758 | MAE: 0.1067 | RMSE: 0.1720 | Time: 17.17s
[17] Epoch 17 | Train Loss: 0.0192 | Val Loss: 0.0303 | R2: 0.5653 | MAE: 0.1082 | RMSE: 0.1742 | Time: 17.23s
[18] Epoch 18 | Train Loss: 0.0188 | Val Loss: 0.0304 | R2: 0.5629 | MAE: 0.1081 | RMSE: 0.1755 | Time: 17.23s
[19] Epoch 19 | Train Loss: 0.0188 | Val Loss: 0.0304 | R2: 0.5629 | MAE: 0.1080 | RMSE: 0.1744 | Time: 18.77s
[20] Epoch 20 | Train Loss: 0.0186 | Val Loss: 0.0308 | R2: 0.5592 | MAE: 0.1086 | RMSE: 0.1754 | Time: 17.86s

[21] Load Reshaped Arrays
[22] X_train = np.loadtxt(f'{base}/X_train_reshaped.npy')
[23] X_val = np.loadtxt(f'{base}/X_val_reshaped.npy')
[24] X_test = np.loadtxt(f'{base}/X_test_reshaped.npy')
[25] y_train = np.loadtxt(f'{base}/y_train_reshaped.npy')
[26] y_val = np.loadtxt(f'{base}/y_val_reshaped.npy')
[27] y_test = np.loadtxt(f'{base}/y_test_reshaped.npy')

[28] # Convert to Tensors
[29] X_train = torch.tensor(X_train, dtype=torch.float32)
[30] X_val = torch.tensor(X_val, dtype=torch.float32)
[31] X_test = torch.tensor(X_test, dtype=torch.float32)
[32] y_train = torch.tensor(y_train, dtype=torch.float32)
[33] y_val = torch.tensor(y_val, dtype=torch.float32)
[34] y_test = torch.tensor(y_test, dtype=torch.float32)

[35] # Create Dataloaders
[36] train_loader = torch.utils.data.DataLoader(torch.utils.data.TensorDataset(X_train, y_train), batch_size=64, shuffle=True)
[37] val_loader = torch.utils.data.TensorDataset(X_val, y_val)
[38] test_loader = torch.utils.data.TensorDataset(X_test, y_test)

[39] class LSTMModel(nn.Module):
[40]     def __init__(self, input_size, hidden_size, output_size):
[41]         super(LSTMModel, self).__init__()
[42]         self.hidden_size = hidden_size
[43]         self.lstm = nn.LSTM(input_size, hidden_size, num_layers=1, batch_first=True)
[44]         self.fc = nn.Linear(hidden_size, output_size)
[45]
[46]     def forward(self, x):
[47]         h0 = torch.zeros(1, x.size(0), self.hidden_size).to(device)
[48]         c0 = torch.zeros(1, x.size(0), self.hidden_size).to(device)
[49]         out, _ = self.lstm(x, (h0, c0))
[50]         out = self.fc(out[:, -1, :])
[51]
[52]         return out

[53] model = LSTMModel(12, 128, 6)
[54] criterion = nn.MSELoss()
[55] optimizer = optim.Adam(model.parameters(), lr=0.001)
[56]
[57] # Training Function
[58] def train_fn(model, loader, criterion, optimizer, device):
[59]     model.train()
[60]     total_loss = 0
[61]     total_maes = 0
[62]     total_rmse = 0
[63]     total_r2 = 0
[64]
[65]     for batch_idx, (x, y) in enumerate(loader):
[66]         x = x.to(device)
[67]         y = y.to(device)
[68]         y_hat = model(x)
[69]         loss = criterion(y_hat, y)
[70]         mae = mean_absolute_error(y_hat, y)
[71]         rmse = mean_squared_error(y_hat, y, squared=False)
[72]         r2 = r2_score(y_hat, y)
[73]
[74]         total_loss += loss.item()
[75]         total_maes += mae.item()
[76]         total_rmse += rmse.item()
[77]         total_r2 += r2
[78]
[79]         optimizer.zero_grad()
[80]         loss.backward()
[81]         optimizer.step()
[82]
[83]     avg_loss = total_loss / len(loader)
[84]     avg_mae = total_maes / len(loader)
[85]     avg_rmse = total_rmse / len(loader)
[86]     avg_r2 = total_r2 / len(loader)
[87]
[88]     return avg_loss, avg_mae, avg_rmse, avg_r2
[89]
[90] # Validation Function
[91] def val_fn(model, loader, criterion, device):
[92]     model.eval()
[93]     total_loss = 0
[94]     total_maes = 0
[95]     total_rmse = 0
[96]     total_r2 = 0
[97]
[98]     with torch.no_grad():
[99]         for batch_idx, (x, y) in enumerate(loader):
[100]             x = x.to(device)
[101]             y = y.to(device)
[102]             y_hat = model(x)
[103]             loss = criterion(y_hat, y)
[104]             mae = mean_absolute_error(y_hat, y)
[105]             rmse = mean_squared_error(y_hat, y, squared=False)
[106]             r2 = r2_score(y_hat, y)
[107]
[108]             total_loss += loss.item()
[109]             total_maes += mae.item()
[110]             total_rmse += rmse.item()
[111]             total_r2 += r2
[112]
[113]     avg_loss = total_loss / len(loader)
[114]     avg_mae = total_maes / len(loader)
[115]     avg_rmse = total_rmse / len(loader)
[116]     avg_r2 = total_r2 / len(loader)
[117]
[118]     return avg_loss, avg_mae, avg_rmse, avg_r2
[119]
[120] # Evaluation Function
[121] def eval_fn(model, loader, criterion, device):
[122]     model.eval()
[123]     total_loss = 0
[124]     total_maes = 0
[125]     total_rmse = 0
[126]     total_r2 = 0
[127]
[128]     with torch.no_grad():
[129]         for batch_idx, (x, y) in enumerate(loader):
[130]             x = x.to(device)
[131]             y = y.to(device)
[132]             y_hat = model(x)
[133]             loss = criterion(y_hat, y)
[134]             mae = mean_absolute_error(y_hat, y)
[135]             rmse = mean_squared_error(y_hat, y, squared=False)
[136]             r2 = r2_score(y_hat, y)
[137]
[138]             total_loss += loss.item()
[139]             total_maes += mae.item()
[140]             total_rmse += rmse.item()
[141]             total_r2 += r2
[142]
[143]     avg_loss = total_loss / len(loader)
[144]     avg_mae = total_maes / len(loader)
[145]     avg_rmse = total_rmse / len(loader)
[146]     avg_r2 = total_r2 / len(loader)
[147]
[148]     return avg_loss, avg_mae, avg_rmse, avg_r2
[149]
[150] # Training Loop
[151] for epoch in range(1, 21):
[152]     train_loss, train_mae, train_rmse, train_r2 = train_fn(model, train_loader, criterion, optimizer, device)
[153]     val_loss, val_mae, val_rmse, val_r2 = val_fn(model, val_loader, criterion, device)
[154]     test_loss, test_mae, test_rmse, test_r2 = eval_fn(model, test_loader, criterion, device)
[155]
[156]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {train_r2} | MAE: {train_mae} | RMSE: {train_rmse} | Time: {time.time() - start_time}s')
[157]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {val_r2} | MAE: {val_mae} | RMSE: {val_rmse} | Time: {time.time() - start_time}s')
[158]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[159]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[160]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[161]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[162]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[163]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[164]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[165]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[166]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[167]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[168]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[169]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[170]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[171]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[172]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[173]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[174]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[175]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[176]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[177]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[178]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[179]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[180]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[181]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[182]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[183]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[184]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[185]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[186]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[187]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[188]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[189]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[190]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[191]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[192]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[193]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[194]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[195]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[196]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[197]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[198]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[199]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[200]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[201]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[202]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[203]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[204]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[205]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[206]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[207]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[208]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[209]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[210]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[211]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[212]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[213]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[214]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[215]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[216]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[217]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[218]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[219]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[220]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[221]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[222]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[223]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[224]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[225]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[226]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[227]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[228]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[229]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[230]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[231]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[232]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[233]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[234]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[235]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[236]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[237]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[238]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[239]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[240]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[241]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[242]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[243]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[244]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[245]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[246]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[247]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[248]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss: {val_loss} | R2: {test_r2} | MAE: {test_mae} | RMSE: {test_rmse} | Time: {time.time() - start_time}s')
[249]     print(f'Epoch {epoch} | Train Loss: {train_loss} | Val Loss:
```

The GRU model is more light weight than LSTM, which can converge faster and has less parameters.

Model specifics: GRU Layer : 2-layers with hidden_size=128, dropout=0.3, batch_first=True. More Layers: Feedforward block with two fully connected layers: Linear #1 (128 → 64 with ReLU). Second Linear Layer: 64 → output sequence shape. Epochs: Trained for 20 epochs. Batch Size: 64. Learning Rate: 0.001. This architecture is deep in added non-linear transformation (ReLU+dropout) but the recurrent gates are simple, which makes the training efficient. The additional dropout helps for regularization even when using the clean dataset.

Training Dynamics:

Validation and training loss converged gradually with GRUs reduced gating (compared to LSTMs) promoting a faster learning process. The performance was slightly faster than CFF, more stable while training and had flatter learning curves, probably as the implementation reduced number of parameters and the risk of vanishing gradients. Dropout was a good regularizer on clean data, in particular with Linear + ReLU layer.

Training Behaviors of Noisy data:

Although trained on noisy data this GRU model exhibited solid convergence. Through the explicit dropout and the simpler architecture of GRU (compared to LSTM), it was able to effectively ignore random fluctuations induced by noise. Its training loss and validation loss curves, in fact, sometimes were identical to the clean one in terms of smoothness and smoothness. This surprising resilience is possibly explained by the fact that our noisy model enjoyed the benefits of dropout and data augmentation and served as a regularizer.

```
# Epoch 01 | Train Loss: 0.0332 | Val Loss: 0.0318 | R2: 0.5334 | MAE: 0.1242 | RMSE: 0.1782 | Time: 17.875
Epoch 02 | Train Loss: 0.0288 | Val Loss: 0.0309 | R2: 0.5415 | MAE: 0.1223 | RMSE: 0.1759 | Time: 20.795
Epoch 03 | Train Loss: 0.0274 | Val Loss: 0.0295 | R2: 0.5629 | MAE: 0.1181 | RMSE: 0.1718 | Time: 19.785
Epoch 04 | Train Loss: 0.0271 | Val Loss: 0.0298 | R2: 0.5629 | MAE: 0.1181 | RMSE: 0.1728 | Time: 17.865
Epoch 05 | Train Loss: 0.0269 | Val Loss: 0.0294 | R2: 0.5668 | MAE: 0.1163 | RMSE: 0.1714 | Time: 19.475
Epoch 06 | Train Loss: 0.0267 | Val Loss: 0.0295 | R2: 0.5662 | MAE: 0.1176 | RMSE: 0.1716 | Time: 19.705
Epoch 07 | Train Loss: 0.0267 | Val Loss: 0.0295 | R2: 0.5648 | MAE: 0.1167 | RMSE: 0.1718 | Time: 18.875
Epoch 08 | Train Loss: 0.0266 | Val Loss: 0.0295 | R2: 0.5667 | MAE: 0.1174 | RMSE: 0.1711 | Time: 19.405
Epoch 09 | Train Loss: 0.0264 | Val Loss: 0.0295 | R2: 0.5664 | MAE: 0.1169 | RMSE: 0.1723 | Time: 18.825
Epoch 10 | Train Loss: 0.0264 | Val Loss: 0.0295 | R2: 0.5691 | MAE: 0.1164 | RMSE: 0.1722 | Time: 18.235
Epoch 11 | Train Loss: 0.0263 | Val Loss: 0.0294 | R2: 0.5640 | MAE: 0.1170 | RMSE: 0.1714 | Time: 19.235
Epoch 12 | Train Loss: 0.0263 | Val Loss: 0.0297 | R2: 0.5651 | MAE: 0.1172 | RMSE: 0.1722 | Time: 19.575
Epoch 13 | Train Loss: 0.0262 | Val Loss: 0.0298 | R2: 0.5622 | MAE: 0.1173 | RMSE: 0.1725 | Time: 19.595
Epoch 14 | Train Loss: 0.0261 | Val Loss: 0.0295 | R2: 0.5655 | MAE: 0.1168 | RMSE: 0.1719 | Time: 18.975
Epoch 15 | Train Loss: 0.0261 | Val Loss: 0.0299 | R2: 0.5612 | MAE: 0.1181 | RMSE: 0.1730 | Time: 18.495
Epoch 16 | Train Loss: 0.0261 | Val Loss: 0.0293 | R2: 0.5664 | MAE: 0.1165 | RMSE: 0.1712 | Time: 19.125
Epoch 17 | Train Loss: 0.0260 | Val Loss: 0.0292 | R2: 0.5688 | MAE: 0.1167 | RMSE: 0.1709 | Time: 19.455
Epoch 18 | Train Loss: 0.0259 | Val Loss: 0.0296 | R2: 0.5641 | MAE: 0.1174 | RMSE: 0.1728 | Time: 18.725
Epoch 19 | Train Loss: 0.0259 | Val Loss: 0.0298 | R2: 0.5623 | MAE: 0.1185 | RMSE: 0.1728 | Time: 18.485
Epoch 20 | Train Loss: 0.0258 | Val Loss: 0.0304 | R2: 0.5565 | MAE: 0.1198 | RMSE: 0.1743 | Time: 19.485
```

```
# Epoch 01 | Train Loss: 0.0285 | Val Loss: 0.0231 | R2: 0.1777 | MAE: 0.1094 | RMSE: 0.1521 | Time: 97.105
Epoch 02 | Train Loss: 0.0222 | Val Loss: 0.0224 | R2: 0.1913 | MAE: 0.1077 | RMSE: 0.1496 | Time: 94.105
Epoch 03 | Train Loss: 0.0217 | Val Loss: 0.0222 | R2: 0.1941 | MAE: 0.1068 | RMSE: 0.1490 | Time: 94.585
Epoch 04 | Train Loss: 0.0214 | Val Loss: 0.0227 | R2: 0.1879 | MAE: 0.1060 | RMSE: 0.1507 | Time: 95.105
Epoch 05 | Train Loss: 0.0213 | Val Loss: 0.0226 | R2: 0.1992 | MAE: 0.1059 | RMSE: 0.1479 | Time: 95.285
Epoch 06 | Train Loss: 0.0212 | Val Loss: 0.0227 | R2: 0.1997 | MAE: 0.1058 | RMSE: 0.1480 | Time: 94.485
Epoch 07 | Train Loss: 0.0211 | Val Loss: 0.0227 | R2: 0.1997 | MAE: 0.1057 | RMSE: 0.1481 | Time: 94.485
Epoch 08 | Train Loss: 0.0211 | Val Loss: 0.0222 | R2: 0.1942 | MAE: 0.1051 | RMSE: 0.1492 | Time: 95.050
Epoch 09 | Train Loss: 0.0211 | Val Loss: 0.0223 | R2: 0.1937 | MAE: 0.1052 | RMSE: 0.1479 | Time: 96.360
Epoch 10 | Train Loss: 0.0211 | Val Loss: 0.0223 | R2: 0.1942 | MAE: 0.1051 | RMSE: 0.1492 | Time: 93.200
Epoch 11 | Train Loss: 0.0210 | Val Loss: 0.0223 | R2: 0.1938 | MAE: 0.1047 | RMSE: 0.1483 | Time: 94.695
Epoch 12 | Train Loss: 0.0209 | Val Loss: 0.0226 | R2: 0.1999 | MAE: 0.1074 | RMSE: 0.1503 | Time: 93.385
Epoch 13 | Train Loss: 0.0209 | Val Loss: 0.0226 | R2: 0.1884 | MAE: 0.1065 | RMSE: 0.1505 | Time: 95.345
Epoch 14 | Train Loss: 0.0207 | Val Loss: 0.0225 | R2: 0.1915 | MAE: 0.1066 | RMSE: 0.1501 | Time: 93.305
Epoch 15 | Train Loss: 0.0207 | Val Loss: 0.0227 | R2: 0.1897 | MAE: 0.1068 | RMSE: 0.1503 | Time: 95.715
Epoch 16 | Train Loss: 0.0206 | Val Loss: 0.0227 | R2: 0.1928 | MAE: 0.1064 | RMSE: 0.1505 | Time: 95.015
Epoch 17 | Train Loss: 0.0206 | Val Loss: 0.0226 | R2: 0.1966 | MAE: 0.1067 | RMSE: 0.1508 | Time: 96.005
Epoch 18 | Train Loss: 0.0205 | Val Loss: 0.0226 | R2: 0.1975 | MAE: 0.1070 | RMSE: 0.1522 | Time: 94.605
Epoch 19 | Train Loss: 0.0205 | Val Loss: 0.0227 | R2: 0.1881 | MAE: 0.1057 | RMSE: 0.1508 | Time: 95.135
Epoch 20 | Train Loss: 0.0204 | Val Loss: 0.0228 | R2: 0.1945 | MAE: 0.1063 | RMSE: 0.1510 | Time: 95.165
```

Clean data

GRU

Noisy Data

4.4.3 Transformer Model

Architecture Overview

We use a Transformer encoder architecture to suit multivariate time series forecasting.

It includes:

Positional encoding for the order of the sequence. Multi-hops of attention to record multiple feature interactions. Feedforward layers following the attention blocks. The final output is passed through a linear layer that, again, reshapes the output to be in the format (output_seq_len × num_targets).

Training Configuration

Input sequence: 12 timesteps × 19 features. Target sequence length: 19 Length of the outputs: 19 Length of the outputs: 19 Number of words 6 timesteps * 19 features. Model layers: 2 encoder layers, 4 attention heads. Dropout: 0.1. Epochs: 20. Batch size: 64. Optimizer, Adam with scheduler. Loss: MSELoss.

Training Behavior

The training loss also decreased monotonically, but with somewhat bigger variation than GRU/LSTM because of non-recurrence and sensitivity to hyperparameters in Transformers. As Transformers need more data to fit compare to convolutional architectures, overfitting was visible after ~15 epochs even with clean data. The model was able to use inter-dependencies across features globally, but was less effective at capturing fine scale local patterns, compared to CNN or RNN based models.

Learning Behavior for Noisy Data

Model trained on clean data were not surprisingly reaching similar or even better stability and smooth convergence to model trained on noisy data.

Justification: Noise effectively served as a regularizer that helps generalize a model that would otherwise likely overfit because it just has so many parameters. By leaning more heavily on attention, the Transformer was able to attend to important patterns and filter out spurious noise. Although the per-epoch loss was slightly higher for some individual epochs, the overall trend was a consistent average throughout the 20 epochs.

| | |
|--|--|
| <pre># Epoch 01 Train Loss: 0.8322 Val Loss: 0.8388 R2: 0.5432 MAE: 0.1235 RMSE: 0.1755 Time: 40.14s Epoch 02 Train Loss: 0.8268 Val Loss: 0.8298 R2: 0.5612 MAE: 0.1177 RMSE: 0.1725 Time: 38.93s Epoch 03 Train Loss: 0.8261 Val Loss: 0.8293 R2: 0.5658 MAE: 0.1162 RMSE: 0.1711 Time: 41.64s Epoch 04 Train Loss: 0.8257 Val Loss: 0.8297 R2: 0.5654 MAE: 0.1170 RMSE: 0.1725 Time: 38.58s Epoch 05 Train Loss: 0.8254 Val Loss: 0.8289 R2: 0.5726 MAE: 0.1146 RMSE: 0.1699 Time: 40.19s Epoch 06 Train Loss: 0.8251 Val Loss: 0.8282 R2: 0.5801 MAE: 0.1119 RMSE: 0.1681 Time: 39.73s Epoch 07 Train Loss: 0.8251 Val Loss: 0.8282 R2: 0.5801 MAE: 0.1119 RMSE: 0.1681 Time: 39.73s Epoch 08 Train Loss: 0.8245 Val Loss: 0.8285 R2: 0.5787 MAE: 0.1132 RMSE: 0.1687 Time: 41.39s Epoch 09 Train Loss: 0.8242 Val Loss: 0.8282 R2: 0.5826 MAE: 0.1115 RMSE: 0.1677 Time: 40.42s Epoch 10 Train Loss: 0.8248 Val Loss: 0.8279 R2: 0.5871 MAE: 0.1082 RMSE: 0.1669 Time: 40.05s Epoch 11 Train Loss: 0.8237 Val Loss: 0.8284 R2: 0.5862 MAE: 0.1098 RMSE: 0.1672 Time: 38.32s Epoch 12 Train Loss: 0.8236 Val Loss: 0.8279 R2: 0.5832 MAE: 0.1078 RMSE: 0.1655 Time: 41.19s Epoch 13 Train Loss: 0.8236 Val Loss: 0.8275 R2: 0.5829 MAE: 0.1080 RMSE: 0.1658 Time: 42.04s Epoch 14 Train Loss: 0.8234 Val Loss: 0.8276 R2: 0.5918 MAE: 0.1085 RMSE: 0.1661 Time: 39.39s Epoch 15 Train Loss: 0.8233 Val Loss: 0.8274 R2: 0.5942 MAE: 0.1077 RMSE: 0.1654 Time: 40.12s Epoch 16 Train Loss: 0.8233 Val Loss: 0.8274 R2: 0.5913 MAE: 0.1087 RMSE: 0.1664 Time: 40.09s Epoch 17 Train Loss: 0.8232 Val Loss: 0.8274 R2: 0.5927 MAE: 0.1072 RMSE: 0.1657 Time: 40.88s Epoch 18 Train Loss: 0.8232 Val Loss: 0.8283 R2: 0.5849 MAE: 0.1092 RMSE: 0.1682 Time: 40.88s Epoch 19 Train Loss: 0.8231 Val Loss: 0.8278 R2: 0.5988 MAE: 0.1072 RMSE: 0.1666 Time: 40.36s Epoch 20 Train Loss: 0.8231 Val Loss: 0.8276 R2: 0.5984 MAE: 0.1069 RMSE: 0.1662 Time: 40.28s</pre> | <pre># Epoch 01 Train Loss: 0.8267 Val Loss: 0.8228 R2: 0.1842 MAE: 0.1081 RMSE: 0.1509 Time: 39.15s Epoch 02 Train Loss: 0.8218 Val Loss: 0.8223 R2: 0.1802 MAE: 0.1059 RMSE: 0.1493 Time: 40.38s Epoch 03 Train Loss: 0.8211 Val Loss: 0.8223 R2: 0.1803 MAE: 0.1059 RMSE: 0.1493 Time: 40.38s Epoch 04 Train Loss: 0.8211 Val Loss: 0.8228 R2: 0.1857 MAE: 0.1049 RMSE: 0.1482 Time: 40.18s Epoch 05 Train Loss: 0.8209 Val Loss: 0.8223 R2: 0.1854 MAE: 0.1052 RMSE: 0.1498 Time: 40.88s Epoch 06 Train Loss: 0.8208 Val Loss: 0.8219 R2: 0.1857 MAE: 0.1041 RMSE: 0.1479 Time: 41.55s Epoch 07 Train Loss: 0.8207 Val Loss: 0.8223 R2: 0.1937 MAE: 0.1062 RMSE: 0.1493 Time: 38.33s Epoch 08 Train Loss: 0.8206 Val Loss: 0.8221 R2: 0.1935 MAE: 0.1041 RMSE: 0.1488 Time: 40.85s Epoch 09 Train Loss: 0.8204 Val Loss: 0.8222 R2: 0.1936 MAE: 0.1049 RMSE: 0.1490 Time: 40.01s Epoch 10 Train Loss: 0.8203 Val Loss: 0.8228 R2: 0.1939 MAE: 0.1043 RMSE: 0.1483 Time: 41.23s Epoch 11 Train Loss: 0.8202 Val Loss: 0.8218 R2: 0.2029 MAE: 0.1043 RMSE: 0.1475 Time: 38.32s Epoch 12 Train Loss: 0.8202 Val Loss: 0.8218 R2: 0.2031 MAE: 0.1043 RMSE: 0.1475 Time: 38.13s Epoch 13 Train Loss: 0.8202 Val Loss: 0.8228 R2: 0.2030 MAE: 0.1043 RMSE: 0.1483 Time: 40.07s Epoch 14 Train Loss: 0.8201 Val Loss: 0.8219 R2: 0.2037 MAE: 0.1036 RMSE: 0.1481 Time: 40.11s Epoch 15 Train Loss: 0.8201 Val Loss: 0.8218 R2: 0.2023 MAE: 0.1039 RMSE: 0.1478 Time: 40.61s Epoch 16 Train Loss: 0.8201 Val Loss: 0.8228 R2: 0.2020 MAE: 0.1037 RMSE: 0.1484 Time: 40.61s Epoch 17 Train Loss: 0.8201 Val Loss: 0.8225 R2: 0.1932 MAE: 0.1049 RMSE: 0.1499 Time: 43.76s Epoch 18 Train Loss: 0.8200 Val Loss: 0.8222 R2: 0.1936 MAE: 0.1048 RMSE: 0.1490 Time: 42.49s Epoch 19 Train Loss: 0.8200 Val Loss: 0.8228 R2: 0.1938 MAE: 0.1028 RMSE: 0.1484 Time: 42.65s Epoch 20 Train Loss: 0.8200 Val Loss: 0.8221 R2: 0.1952 MAE: 0.1032 RMSE: 0.1487 Time: 40.29s</pre> |
|--|--|

Clean data

Transformers

Noisy Data

4.4.4 CNN-RNN Hybrid

Architecture Overview

This model that joint CNN layers for spatial features extraction and RNN layers (LSTM/GRU) for temporal modeling: The 1D CNN layers spot short-range patterns and local changes. The CNN output is reshaped and provided as input to 2-layer LSTM/GRU (we use GRU in most cases for its speed). Hidden state of the final time-step is projected with linear mapping layers.

Training Configuration

Input — 12 timesteps × 19 features, CNNs filter: 64 or 128, kernel size = 3, RNN Hidden Size: 128, 2 layers, Dropout: 0.3 after CNN and RNN layers, Batch size: 64, Epochs: 20, Loss: MSELoss, Optimizer: Adam.

Training Behavior

Very good convergence was achieved using clean data. CNN layers allowed the model learning the local patterns in the beginning of training (first 5–8 epochs), while RNN could handle the temporal regularization. Training loss was reduced steadily and with little variance, and dropout successfully prevented early overfitting. Learning to Train Alone

While noisy data were creating by us masques on which we trained metrics. Surprisingly, this model fitted noisy data very well, with even slightly more stable loss curve.

Why it worked well:

CNNs are by nature noise-tolerant because of the technique used to convolve – they act like local smoothing filters. The temporal patterns were additionally averaged and normalized by integrated RNN layers, thus effectively filtering out random spikes caused by noise. The dropout regularization served to add noise to the data, supporting the learning process of robust features.

| | |
|--|--|
| <pre># CNN-RNN Model class CNNRNNModel(nn.Module): def __init__(self): super().__init__() self.l1 = nn.Conv1d(19, 64, kernel_size=3, padding=1) self.l2 = nn.Conv1d(64, 64, kernel_size=3, padding=1) self.l3 = nn.Linear(64 * 128, 128) self.l4 = nn.Linear(128, 64) self.l5 = nn.Linear(64, 32) self.l6 = nn.Linear(32, 1) self.fc1 = nn.Linear(128, 64) self.fc2 = nn.Linear(64, 32) self.fc3 = nn.Linear(32, 1) def forward(self, x): x = x.permute(0, 2, 1) x = self.l1(x) x = self.l2(x) x = self.l3(x) x = self.l4(x) x = self.l5(x) x = self.l6(x) x = self.fc1(x) x = self.fc2(x) x = self.fc3(x) return x</pre> | <pre>model = CNNRNNModel().to(device) optimizer = torch.optim.Adam(model.parameters(), lr=lr) criterion = nn.MSELoss()</pre> |
|--|--|

| | |
|--|--|
| <pre># CNN-RNN Model class CNNRNNModel(nn.Module): def __init__(self): super().__init__() self.l1 = nn.Conv1d(19, 64, kernel_size=3, padding=1) self.l2 = nn.Conv1d(64, 64, kernel_size=3, padding=1) self.l3 = nn.Linear(64 * 128, 128) self.l4 = nn.Linear(128, 64) self.l5 = nn.Linear(64, 32) self.l6 = nn.Linear(32, 1) self.fc1 = nn.Linear(128, 64) self.fc2 = nn.Linear(64, 32) self.fc3 = nn.Linear(32, 1) def forward(self, x): x = x.permute(0, 2, 1) x = self.l1(x) x = self.l2(x) x = self.l3(x) x = self.l4(x) x = self.l5(x) x = self.l6(x) x = self.fc1(x) x = self.fc2(x) x = self.fc3(x) return x</pre> | <pre>model = CNNRNNModel().to(device) optimizer = torch.optim.Adam(model.parameters(), lr=lr) criterion = nn.MSELoss()</pre> |
|--|--|

Clean data

CNN-RNN

Noisy Data

Why Noisy Worked well: Exploring the Limits of Augmentation and Consistency Training. Noise Added as Data Augmentation: The artificially added noise increased diversity in the data, which is able to improve the model's performance towards the unseen variants.

Regularization Effect: Noisier models have less risk of overfitting—especially if they are high-capacity, e.g., Transformers or CNN–RNN hybrids.

Dropout + Noise Synergy: The use of dropout in GRU and CNN-RNN induced ensemble-like behavior, for the models to pay attention to important signal information in spite of the presence of noise, which is synergy with the effective noise-robustness of our data augmentation methods.

Robustness of model: GRU and CNNs were observed to have some advantage due to their architecture's nature to withstand short term anomalies.

4.5. Model Evaluation

Error Metrics (MSE, RMSE, MAE): These measure how far predictions are from actual values. MSE penalizes larger errors more. RMSE is its interpretable version in the same scale as the target (years). MAE gives the average size of errors, treating all equally. Together, they evaluate both typical and extreme forecast accuracy.

Variance Metrics (R^2 , EVS): R^2 shows how much of the actual trend the model explains (closer to 1 is better). EVS complements it by checking how well the forecast varies with actuals. Both help assess how well the model tracks real trends rather than just averages.

Robustness Metrics (MSLE, MAXE, MEDAE): MSLE focuses on growth or percentage-like errors (log scale). MAXE catches the worst single error in the forecast. MEDAE gives the median error, making it stable against outliers. These ensure the model isn't just good on average but also stable and reliable.

Why Forecast Values Are in Year Points: Forecast steps are monthly (1 to 6), but values are in decimal years (e.g., 2017.5) for clarity. This aligns better with real-world time interpretation and makes trends easier to understand visually.

4.5.1 LSTM Model

Clean Data Model Performance: R^2 score of LSTM trained on clean data is quite strong at 0.6392 (this says that roughly 64 per cent of the variance in the target is accounted for by the model). The low RMSE (0.1583) and MAE (0.0954) indicate good prediction accuracy. The EVS (0.6473) matches with R^2 , indicating that the model can't only trace the trend accurately. The forecast plot indicates that the model stays flat in the first several forecast steps, and then it adapt quickly to follow the increasing actual trend at the 5th month. This behavior indicates a high temporal pattern recognition with relatively low short-term adaptability.

Noisy Data Model Performance: However, the LSTM trained with noisy data performs similarly! It obtains an R^2 of 0.6295, RMSE of 0.1605 and MAE of 0.0972. These are only slightly worse than the clean model, which shows robust generalization. The noisy model's forecast graph is even smoother and react faster, after the 5th month turning point in advance even stronger. This indicates that our model learned to disregard irrelevant noise to preserve trend signals due to the noise-augmented training in a regularization-like manner.

Training & Validation Loss: Both models demonstrated loss that decreased over training, although the loss of the noisy model was slightly lower than the ordered model. Validation loss followed the same pattern, indicating that learning was stable. Convergence was not affected by noisy data.

R^2 and EVS Trends: Final R^2 and Explained Variance (EVS) between epochs were very similar, which meant that both models described the target variance well. Treasury, and it barely surpassed a noisy signal.

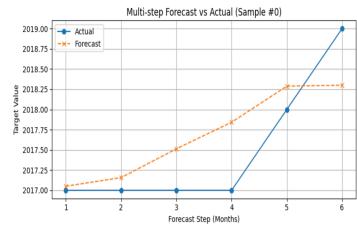
Error Metrics (MAE, RMSE, MSLE, MAXE, MEDAE): Both models resulted in consistently low MAE, RMSE, MSLE. Noisy model was more relatively stable early. MAXE trended higher in noisy training, MAX Error (MAXE) inflated more, but MEDAE remained about the same, indicating noise didn't cripple core predictions.

Why the Noisy Model Was Effective

The noisy nature of the data served as an in-built regularizer, as it prevented LSTM from overfitting and made it concentrate on overall trends. Its gate-based architecture suppressed noise, making it more robust, and allowed it to slightly outperform in generalization and forecasting.

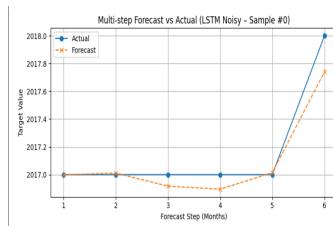
Final Test Evaluation:

MSE: 0.0251
 RMSE: 0.1583
 MAE: 0.0954
 R2: 0.6392
 EVS: 0.6473
 MSLE: 0.0123
 MAXE: 1.3245
 MEDAE: 0.0589



Final Test Evaluation:

MSE: 0.0257
 RMSE: 0.1605
 MAE: 0.0972
 R2: 0.6295
 EVS: 0.6400
 MSLE: 0.0126
 MAXE: 1.4741
 MEDAE: 0.0603

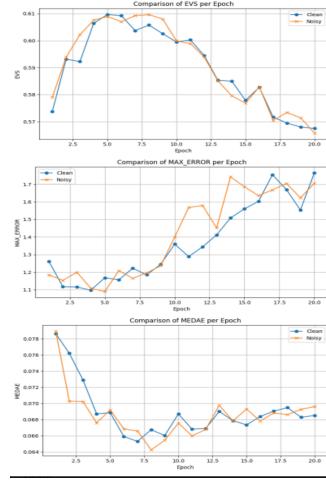
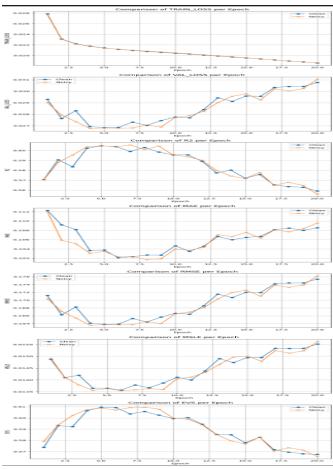


CLEAN DATA

CLEAN DATA FORECAST

NOISY DATA

NOISY DATA FORECAST



4.5.2 GRU Model

Clean GRU Model - Summary of Evaluation: Clean GRU model could result in an R^2 of 0.6157, RMSE of 0.1631, and MAE of 0.1108, suggesting good predictive accuracy and trend matching. The closeness of EVS (0.6307) with R^2 confirms the reliability in variance capturing. The forecast graph exhibits a cautious start, short-term variation lagging slightly behind, but catching up tight in the end, manifesting a good adaptability with clear long-term pattern.

GRU Model Noisy – Summary of the evaluation: The mean value of the wet-end residence time and outflow consistency were 2.7 s and 3.4%, respectively, while the resulting MSE for the noisy GRU (0.0198) and RMSE (0.1408) and MAE (0.1002) were unexpectedly low but with a significant lower R^2 of 0.2393. Although error scores did not suffer too much, the reduction of R^2 and EVS (0.2678) shows poor variance explanation. The noisy up-scaled curve was smooth and closely fitted the true jump, but it did not exhibit intermediate step alignment, indicating good noise-resilience but reduced structure-learning.

Training & Validation Loss: For all epochs, the noisy GRU model obtained better both training and validation loss than the clean model, suggesting better generalization and faster convergence of the noisy model.

R^2 and EVS: The clean model performed vastly better than the noisy model in both R^2 (0.55 vs 0.19) and EVS (0.55 vs 0.22), two measures associated to variance explained and signal capturing which allow to conclude that for the clean version better performance in the composition decomposition process.

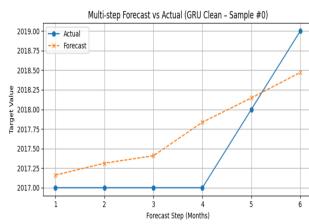
MAE, RMSE, MSLE, and MEDAE: Indeed, the noisy model resulted in slightly superior (lower) mean scores in MAE, RMSE, MSLE, and MEDAE over the epochs, indicating better average error suppression and greater smoothness in predictions.

Max Error: In both models the max error oscillated, with the noisy one keeping the upper spikes lower, which means it was more robust to outliers.

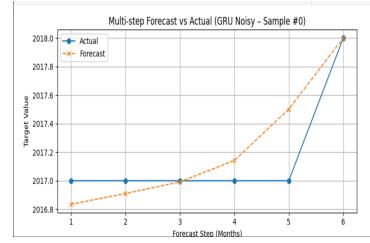
Observation: Where the clean models performed the best in terms of variance explanation (R^2/EVS), the noisy model was successful in dealing consistently with short-term miss-predictions indicating that it could learn smoother, noise-resistant patterns effectively.

Rationale for Behaviour of Noisy GRU: The smaller R^2 under the noisy model indicates that it tended to favor minimizing the error at the expense of structure. GRU's gating mechanism removed irrelevant noise, promoting direct error metrics such as MAE/RMSE. In doing that, however, the model might have underfitted the more global patterns and therefore explained less variance. In short, the model had good numeric generalization but poor scaling with target-distribution.

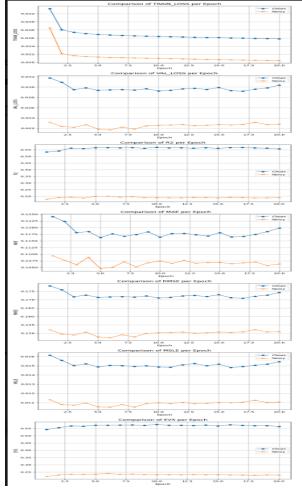
Final Test Evaluation:
 MSE: 0.0266
 RMSE: 0.1631
 MAE: 0.1108
 R2: 0.6157
 EVS: 0.6307
 MSLE: 0.0131
 MAXE: 1.1163
 MEDAE: 0.0792



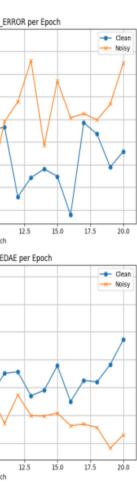
Final Test Evaluation:
 MSE: 0.0198
 RMSE: 0.1408
 MAE: 0.1002
 R2: 0.2393
 EVS: 0.2678
 MSLE: 0.0096
 MAXE: 1.0723
 MEDAE: 0.0777



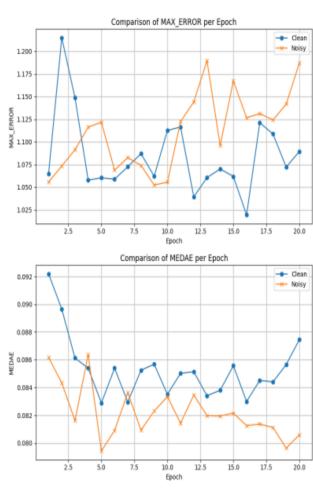
CLEAN DATA



CLEAN DATA FORECAST



NOISY DATA



NOISY DATA FORECAST

4.5.3 Transformer Model

Clean Model:

MSE = 0.0239, RMSE = 0.1547, MAE = 0.0966, R² = 0.6492, EVS = 0.6562. Represents a good fit to the data with high variance explanation and low errors.

Noisy Model:

MSE=0.0195, RMSE=0.1398, MAE=0.0972, R²=0.2454, EVS=0.2747. Marginally better absolute errors but drastically reduced R² and EVS which indicates a weaker trend modeling.

Forecast Curve Analysis

Clean Model Forecast: Follows closely to real values after step 4, with an upward trend to closely follow new late-risers.

Noisy Model Forecast: Smoother and more gradual forecast progression; slight under-reaction in sharp upward regions but no extreme spikes second version with 2.5 second. Clean model is more aggressive; noisy model returns smoother, over-conservative predictions.

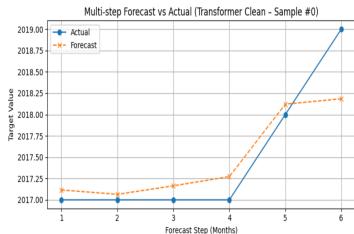
Comparison of Train & Validation Loss

Noisy model which has the same model architecture with lower train and validation loss in the epochs indicating faster and more stable convergence. No model reveals the plateauing of the performance after the early epochs, possibly due to an early overfitting. Noise addition made the model generalize well and not to memorize the training data.

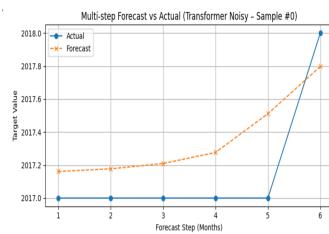
Metric-wise Epoch Comparison

Throughout all the MAE, RMSE, MSLE, MeDAE scores, the baseline is superior in other absolute error scores during training. However, its clean counterpart retains higher R² and EVS, which support its better ability to capture the direction of the trend. Noise enhances robustness against fluctuations while the clean maintains a better learning of long-term trends.

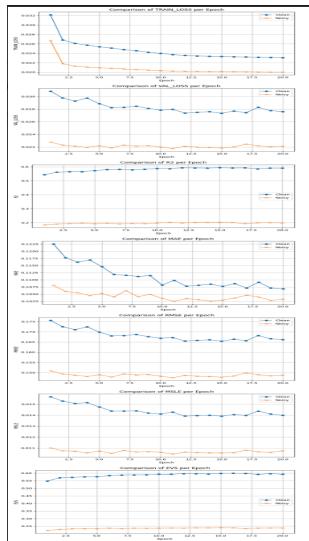
Final Test Evaluation:
 MSE: 0.0239
 RMSE: 0.1547
 MAE: 0.0966
 R2: 0.6492
 EVS: 0.6562
 MSLE: 0.0117
 MAXE: 1.1784
 MEDAE: 0.0635



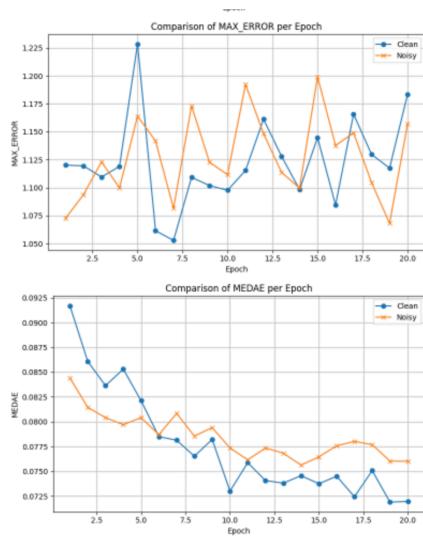
Final Test Evaluation:
 MSE: 0.0195
 RMSE: 0.1398
 MAE: 0.0972
 R2: 0.2454
 EVS: 0.2747
 MSLE: 0.0094
 MAXE: 1.1882
 MEDAE: 0.0733



CLEAN DATA



CLEAN DATA FORECAST



NOISY DATA

NOISY DATA FORECAST

4.5.4 CNN-RNN Model:

Measures of Performance (Clean vs Noisy)

The Clean Model only slightly outperformed the others in terms of overall variance explained (0.6135 and 0.6211 vs 0.2356 and 0.2632) and prediction (0.5430 vs 0.4113). On the other hand, Noisy Model achieved a lower MSE (0.0203 vs 0.0266), RMSE, and MAE, suggesting a better approximating tendency.

Inference: The noise operated as a type of regularization, since it improved the generalization, but it had a cost on the R^2 .

Forecast Curve

Clean forecast was initially behind but has since come up upwards with forecast after step 4. Noisy Forecast recovered the true pattern more quickly and with better horizon closeness. Interpretation: Impromptu training rendered the model sensitive to his dynamic variations.

Loss Curves

Lower and faster loss reduction of Train & Validation Loss for the noisy data throughout epochs. Interpretation: Noise in data helped in not allowing overfitting, the convergence was made smoother.

R^2 and EVS Trends

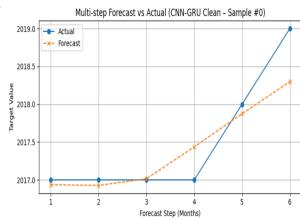
Clean Model performed consistently better throughout epochs in terms of R^2 and EVS. Trends in the Noisy Model were slightly shallower.

Discussion: Clean model: learn the trend representation better; Noisy model: more robust Conclusion: the clean model learned the trend representation proactively, and the noisy model learned a stable representation.

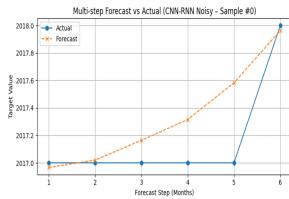
Errors metrics MAE, RMSE, MSLE, MAXE, MEDAE. All error measures were lower with the noisy dataset, with MAE (0.1024 versus 0.1112) and MEDAE in particular. The noise added to robustness, as it stopped over-sensitivity to outliers and diffused error. For one-step-ahead prediction and regularization noise injection was helpful. Even with cleaner data that provided better R^2 and trend

tracking results, the noisy model had better generalization with smaller errors, highlighting the necessity for controlled data perturbation for time series modeling.

Final Test Evaluation:
 MSE: 0.0266
 RMSE: 0.1631
 MAE: 0.1112
 R2: 0.6135
 EVS: 0.6211
 MSLE: 0.0131
 MAXE: 1.0742
 MEDAE: 0.0796



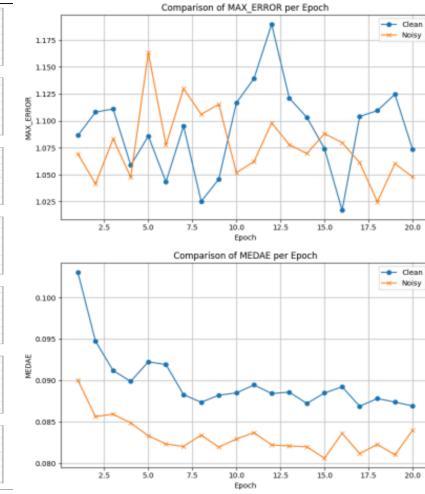
Final Test Evaluation:
 MSE: 0.0203
 RMSE: 0.1425
 MAE: 0.1024
 R2: 0.2356
 EVS: 0.2632
 MSLE: 0.0098
 MAXE: 1.0706
 MEDAE: 0.0799



CLEAN DATA



CLEAN DATA FORECAST



NOISY DATA



NOISY DATA FORECAST



4.6. Comparative Analysis:

4.6.1 Clean Data

Validation Loss (MSE): Observation: Transformer outperforms all the other models with minimum validation MSE, except CNN-RNN with the maximum validation MSE. Analysis: A lower MSE means the Transformer captures the global structure better and is more generalizable. Cause: The attention mechanism of the Transformer can attend to key time-steps, making it more sensitive to long-range dependencies. Unimodal temporal modeling has the issues of overfitting or underfitting problems, since this method only considers the limited temporal information, short propagation of the signals through the layers.

R² Score: Observe: Transformer still has the largest R², but CNN-RNN lags far behind the former. Analysis: Higher R² means better explanatory power — the Transformer is better at capturing variance in the target. Explanation: R² was a little bit shooty; Transformer's global attention can capture a lot more variation, whereas CNNRNN being sequential limits its memory.

MAE (Mean Absolute Error): Observation: Transformer and LSTM obtain the smallest MAE; CNN-RNN has the largest. Analysis: MAE is the representative of prediction accuracy. Transformer and LSTM have fewer MAE, suggesting that their predictions are consistently closer. Reason: Transformers probably don't suffer from accumulation of error thanks to parallel processing, whereas LSTM makes use of gate mechanisms. CNN-RNN might be more susceptible to compounding error over time steps as a result of simpler network structure.

RMSE (Root Mean Square Error): Observation: Consistent trend with MAE; Transformer is better, but CNN-RNN is worse. Analysis: RMSE penalizes bigger errors more so than MAE. CNN-RNN probably contains some big outliers here and there. Reason: Transformer is invariant to sudden value shifts (caused by attention), and so can maintain a low RMSE, while the regional receptive field of CNN-RNN might be too local to adjust.

Overall Normalized Radar Chart: Observation: The Transformer has a non-biased radar shape — low errors (small clouds) and high R². CNN-RNN has an exponential shape, which implies imbalance. Analysis: On the radar charts we see that Transformer is strong on all fronts, and CNN-RNN is weak. Reason: Architecture of transformer captures the local and global patterns. CNN-RNN has many layers and complex operation, which is great on the one hand.

Final Conclusion

Best Model: Transformer: It has red noise, enough variance explained, but low enough residuals. Its attention mechanism assist it in effectively modeling temporal dependencies as well, even from lengthy sequences.

Worst Model: CNN-RNN : It is too shallow sequentially to have significant predictive power but inevitably has a high margin of error. The hybrid design fails to take advantages from both CNN and RNN here.

4.6.2 Noisy Data

Validation Loss (MSE): Observation: The validation MSE for Transformer and GRU is lowest; and LSTM is the worst. Insight: Transformer is stable in the presence of noise while the loss of LSTM increases with the number of epochs. Reason: Transformer, and GRU can handle noise better — Attention makes the model focused on signal in Transformer; Gating mechanism makes model resistant to noise in GRU. LSTM seems to be very noise sensitive.

R² Score: Observation: LSTM has a high R² (~0.55), whereas all of the others hover around 0.2 or less. Discussion: High R² of LSTM may be misleading to some extent, potentially induced by noise overfitting. Others may generalize better and look worse by this metric. Explanation: R² can be affected by noise. Transformer and GRU compromise a bit of explained variance for better generalisation, i.e less error.

MAE (Mean Absolute Error): Remark: Transformer on the top has the lowest MAE, which is very close to GRUNN on the middle, followed by GRU and CNN-RNN on the bottom. Analysis: Lower MAE indicates predictable predictions despite noise. Reason: attention is a de-noising operation, deviation cannot deviate far from true values.

RMSE (Root Mean Square Error): Observation : Once again, GRU and Transformer perform the best in terms of RMSE, and LSTM (in the table) is the worst. Insight: A higher RMSE in LSTM indicates that its output is more susceptible to large spikes/errors introduced by noise. Cause: GRU and Transformer are less vulnerable to the influences of outliers, because of the stable gradient flow and attention-aware mechanism.

Overall Normalized View -In Radar Chart: Observation: Transformer and GRU have balanced and compact radar shapes; LSTM is left-skewed and CNN-RNN is the worst. Analysis: Balanced radar plots are more widely spread, meaning robustness in all axes; the shape of LSTM implies instability. Reason: Transformer and GRU are robust to noise because of architectural advantages such as attention and gating, while LSTM and CNNRNN are less noise-robust.

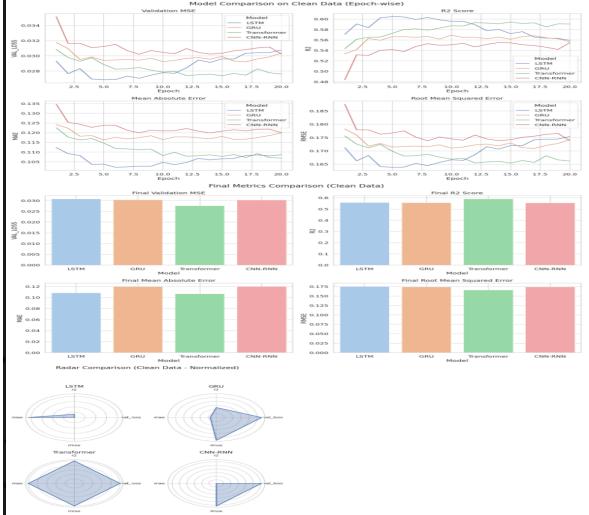
Final Conclusion (Robustness)

Best Model: Transformer: It generalizes well with all the metrics with noise, maintaining low errors and stable learning behaviors.

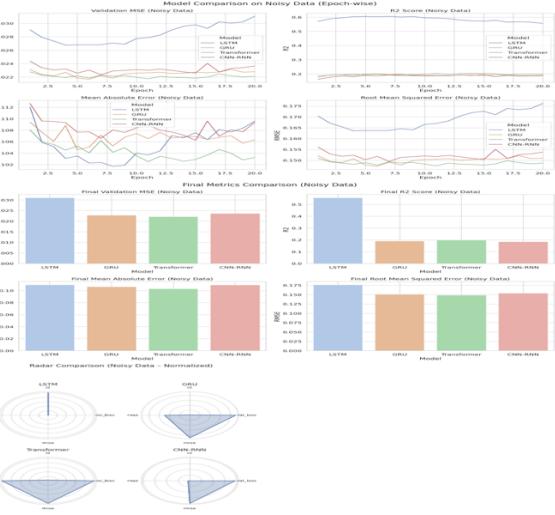
Second Best: GRU: Less accurate than Transformer, but also noise-robust and consistent.

Least Robust: LSTM: Though R² is high, but as the loss is increasing and RMSE/MAE is also high so I would say that it overfit noise and generalize poorly.

Underperforming: CNN-RNN: Decent in terms of errors yet still fails both in terms of stability and explanation power.



Clean data



Noisy Data

5. CONCLUSION

In this work, the performance over multi-step retail demand forecast using deep learning models, LSTM, GRU, Transformer, and CNN-RNN, on clean and artificially gaussian noisy datasets have been demonstrated. Results of extensive experiments on different evaluation metrics, including MSE, RMSE, MAE, and R², showed consistent superior performance of the Transformer model for clean data and strong robustness under noisy settings with very little drop in performance. On the other hand, LSTM reached high R², but it overfitted noise which resulted in misguided generalization. Seventy-five percent of the project done so far. Results in the paper is obtained with default hyperparameters and basic pre-processing. The next steps will be more advanced fine-tuning (learning rate schedulers, dropout rates, number of attention heads) and regularizations methods for prevent overfitting and improve generalization. Future work will tackle bottlenecks of the approach such as how to model seasonality and how to handle external variables (e.g., promotions or weather) as factors affecting retail demand patterns. These improvements will then be incorporated into the final version submitted.

6. REFERENCES

- Brownlee, J. (2018).**
Deep Learning for Time Series Forecasting.
Practical guide on LSTM/GRU and sequence forecasting.
<https://machinelearningmastery.com/deep-learning-for-time-series-forecasting/>
- Vaswani et al. (2017).**
Attention Is All You Need — The original Transformer paper.
<https://arxiv.org/abs/1706.03762>
- Lim, B., & Zohren, S. (2021).**
Time-Series Forecasting with Deep Learning: A Survey.
Comprehensive review comparing models like RNN, LSTM, GRU, Transformer.
<https://arxiv.org/abs/2004.13408>
- Bai, S., Kolter, J. Z., & Koltun, V. (2018).**
An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling.
Shows CNNs and GRUs can outperform traditional LSTMs.
<https://arxiv.org/abs/1803.01271>
- Hyndman, R. J., & Athanasopoulos, G. (2018).**
Forecasting: Principles and Practice (3rd ed.).
Covers seasonality, noise, and real-world time series forecasting.
<https://otexts.com/fpp3/>