

Table 1: Explored Models & Hyperparameters

Model Type	Key Hyperparameters
Linear Regression	n_jobs=7
Adaboost	n_estimators=5,learning_rate=0.1,loss='square'
XGBoost Regressor	n_estimators=100, max_depth=6, learning_rate=0.1, subsample=0.8, colsample_bytree=0.8, tree_method='hist'
LightGBM	params = {"objective": "tweedie", "metric": "rmse", "tweedie_variance_power": trial.suggest_float("tweedie_variance_power", 1.05, 1.3), "learning_rate": trial.suggest_float("learning_rate", 0.01, 0.05), "num_leaves": trial.suggest_int("num_leaves", 128, 512), "min_data_in_leaf": trial.suggest_int("min_data_in_leaf", 256, 2047), "feature_fraction": trial.suggest_float("feature_fraction", 0.6, 0.95), "bagging_fraction": trial.suggest_float("bagging_fraction", 0.6, 0.95), "bagging_freq": 1, "lambda_l2": trial.suggest_float("lambda_l2", 0.0, 1.0), "boost_from_average": True, "max_bin": 127, "bin_construct_sample_cnt": 20000000, "force_row_wise": True, "verbosity": 1, "nthread": 8, "seed": 42}
LSTM (Encoder)	<b>Embedding + LSTM</b> structure for categorical inputs <b>LSTM</b> for sequential numerical features <b>Dense layers with varied activation functions</b> <b>Batch normalization</b> for stabilization <b>Activations:</b> Mix of relu, sigmoid, tanh across layers
CNN LSTM	<b>Embedding dimensions</b> range from 10 to 150 <b>Conv1D filters</b> range from 3 to 7; <b>kernel size</b> = 8 (categorical), 16 (numerical) <b>LSTM units:</b> 50 (Item_Id), 10 (others), 32 (numerical) <b>Dense layers:</b> 256 → 128 → 64 → 1 <b>Batch normalization</b> for stabilization <b>Activations:</b> Mix of relu, sigmoid, tanh across layers

Table 2: Engineered Features & Rationale

Type	Feature	Description	Column_Name	Why do we use	Used in Final Model
Categorical Encoding	item_id	Item ID	item_id	Encoding transforms categorical features into <b>model-readable, numerical formats.</b>	Yes
	dept_id	Department ID	dept_id		Yes
	cat_id	Category ID	cat_id		Yes
	store_id	Store ID	store_id		Yes
	state_id	State ID	state_id		Yes
	event_name_1	Event Name	event_name_1		Yes
	event_name_2	Event Name	event_name_2		Yes
	event_type_1	Event Type	event_type_1		Yes
	event_type_2	Event Type	event_type_2		Yes
	year	Year of sale	year		Yes
Calendar Features	dayofweek	Day of week (1=Saturday)	wday	Helps model seasonality, weekly and monthly cycles.	Yes
	weekofyear	Week number (ISO Week Number?)	weekofyear		Yes
	month	Month number	month		
	year	Year number	year		Yes
	day	Day of month	day		
	weekend	Flag: is it the weekend?	weekend		
	Is month start / end	If month start or end	is_month_start/end		Yes
	Is quarter start / end	If quarter start or end	is_quarter_start/end		Yes
	Is Year start/end	If year start or end	year_start/year_end		Yes
	Season	Season (Winter, Spring, etc)	season		Yes
	Day of year	Day of the year	day_of_year		

Event/Holiday Features	is_event	Flag: Is there an event?	is_event	Strong lift around holidays and events.	
	Event proximity	Days until / since next event	days_since_event/days_until_event		
Promotion/Snap Features	snap_CA, snap_TX, snap_WI	SNAP active flags	snap_active	<b>SNAP (Supplemental Nutrition Assistance Program)</b> is a U.S. government program that helps low-income individuals and families buy food.	Yes
	Days since SNAP active	For SNAP customers, since last promotion active	days_since_snap		
	Days until next SNAP	For SNAP customers, until next promotion active	days_until_next_snap		
Lag Features	sales_lag_1	Yesterday's sales	sales_lag_1	Time-series models must have lag features to capture autocorrelation.	
	sales_lag_7	Last week's sales	sales_lag_7		
	sales_lag_14	Two weeks	sales_lag_14		
	<u>Sales_lag_28_35_42_49_56_63_70_77_84_91_98</u>	Custom lag	direct_lag_n		Yes
Rolling Window Features	rolling_mean_7/14	Moving average	rolling_mean_7/rolling_mean_14	Capture trends and seasonality patterns dynamically.	
	rolling_std_7/14	Rolling standard deviation (volatility)	rolling_std_7/rolling_std_14		
	expanding_mean	Expanding mean for trend, It's the <b>cumulative average</b> of sales up to that point in time (excluding today), calculated per item-store (id).	expanding_mean		
	rolling_min / max	Rolling min/max for local peak detection. It tells you the <b>lowest sales value</b> in the <b>last 7 days</b> before a given day.	rolling_min_7/rolling_max_7		
	Rolling statistical mean of 7, 14, 30, 60, 360 days with shift of 28 days	<b>Rolling statistical features</b> (mean) over past sales data, with a <b>shift of 28 days</b> to <b>prevent data leakage</b> .	roll_n_shift_28_mean		Yes
	Rolling statistical standard deviation of 7, 14, 30, 60, 360 days with shift of 28 days	<b>Rolling statistical features</b> (standard deviation) over past sales data, with a <b>shift of 28 days</b> to <b>prevent data leakage</b> .	roll_n_shift_28_std		Yes
Cumulative /Expanding Features	cumulative_sales	Cumulative sum of sales	cumulative_sales	Good for lifecycle analysis and stockout behavior.	
	cumulative_mean_sales	Expanding mean of sales	cumulative_mean_sales		

	days_since_last_sale	Counter since last sale event (in case of NA counted days since start)	days_since_last_sale		
Hierarchical Aggregated Features	sales at dept/store level	Aggregated sales at dept or store level	dept_sales/store_sales	Aggregations help capture group-level seasonality and shared patterns.	
	sales at state/category level	Aggregated sales at state or category level	state_sales/cat_sales		
	rolling mean at group level	Rolling means for groupings (store-dept)	rolling_mean_store_dept_7/ rolling_mean_store_dept_7_lag_7		
Lag on Engineered Features	Lag of rolling mean	Lag of moving averages	rolling_mean_7_lag_7/ rolling_mean_14_lag_14	Captures Temporal Patterns in Aggregates.	
	Rolling Mean Change	Derive Growth Rates	rolling_mean_change		

Table 3: Model Performance Comparison

Model	Private WRASSE
Linear Regression	0.883
Adaboost	1.094
XGBoost	2.12
LightGBM	0.611
LSTM (Encoder)	0.980
CNN LSTM	1.285
Ensemble	0.616

M5 Forecasting - Accuracy | Kaggle

https://www.kaggle.com/competitions/m5-forecasting-accuracy/leaderboard

Import bookmarks... College Applications Entertainment Important Aus Visa Carlson Career Search Books

+

Create

Home

Competitions

Datasets

Models

Code

Discussions

Learn

More

User Rankings

Documentation

Progression

Host a Competition

Educator Resources

Support/Contact

Community Guidelines

View Active Events

Search

M5 Forecasting - Accuracy

Estimate the unit sales of Walmart retail goods

Overview

Data

Code

Models

Discussion

Leaderboard

Rules

Team

Submissions

Raw Data

Refresh

YOUR RECENT SUBMISSION

✓

🕒

lgbm\_submission\_1.csv

Submitted by Pankaj Nandal · Submitted 16 minutes ago

Score: 0.61123

Public score: 0.73513

Jump to your leaderboard position

Search leaderboard

Public

Private

The private leaderboard is calculated with approximately 50% of the test data. This competition has completed. This leaderboard reflects the final standings.