

A low-angle, upward-looking photograph of a modern building's glass facade. The image shows a grid of glass panels held by metal frames, reflecting the sky and surrounding environment. The perspective creates a sense of height and architectural scale.

FINAL PRESENTATION

Team6

Outline

1. Data Overview
2. Data Preprocessing
3. The Models
4. Result

Data Overview

The Features

Numerical

- ❏ On_promotion
- ❏ Oil price

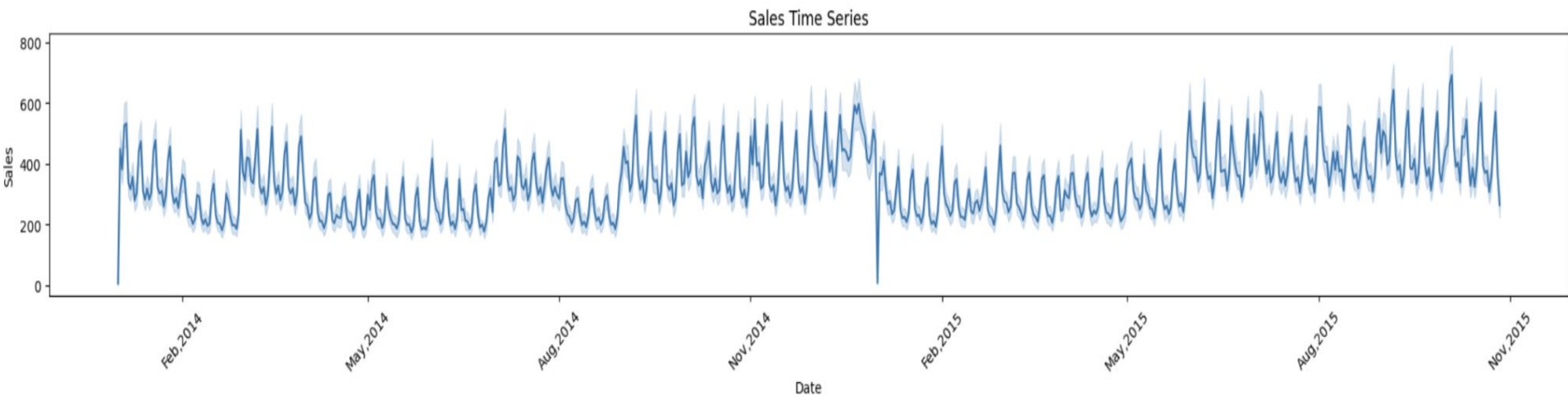
Categorical

- ❏ Store_nbr
- ❏ Family
- ❏ Store type
- ❏ Earthquake
- ❏ Holidays
- ❏ City

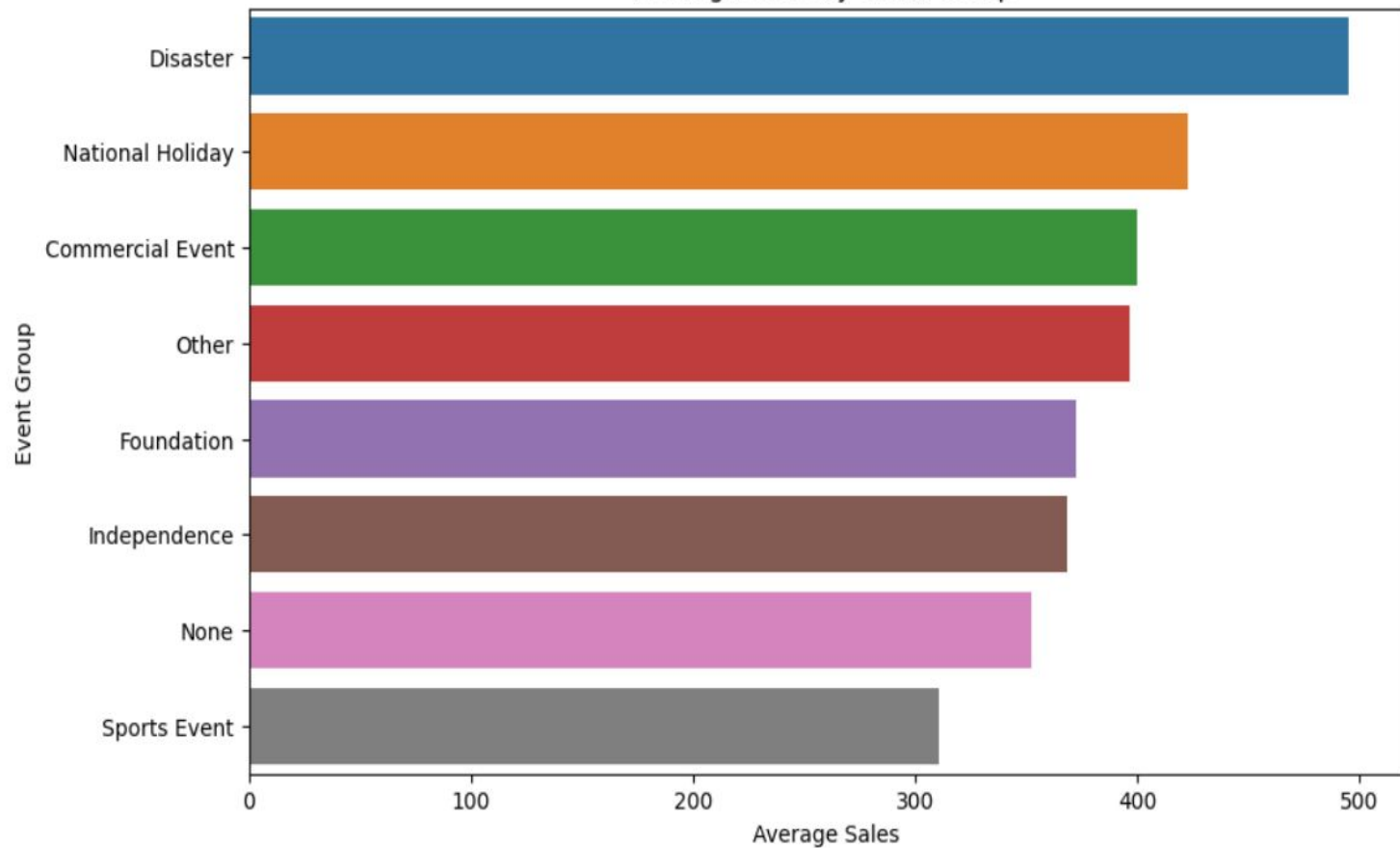
Target Feature



Sales



Average Sales by Event Group



Problem

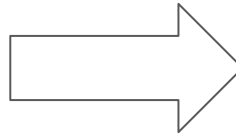
- Predict sales for each store and product category in the retail industry.
- Accurate sales forecasting allows for optimization of inventory management and promotional strategies.

Data Preprocessing and Feature Engineering

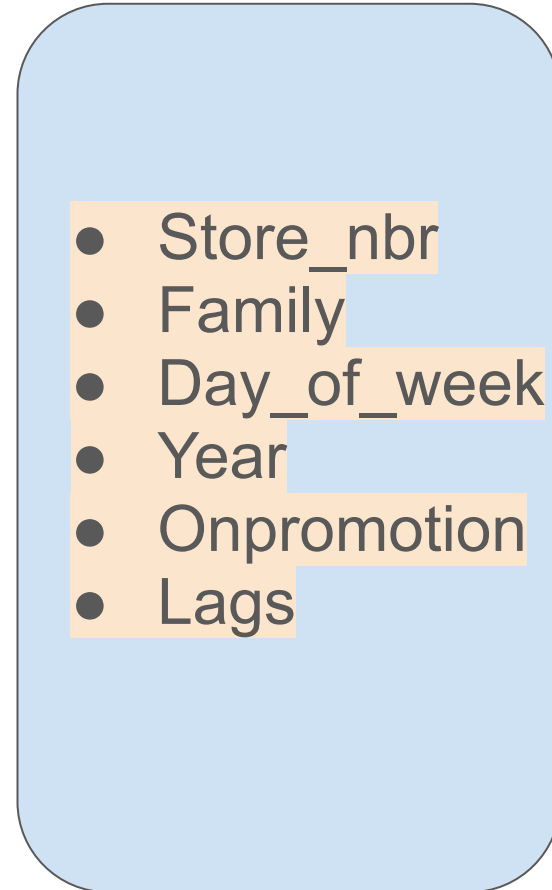
Improvements in Our Approach

- Reduced the number of features
- Created new, effective features

What we used to use:



What we use now:



Specific Improvements

- Removed unnecessary categorical columns
- Used LightGBM for efficient handling of categorical data



Simplified Date Features

- Focused on two features: `day_of_week` and `year`
- Sufficient for effective model learning

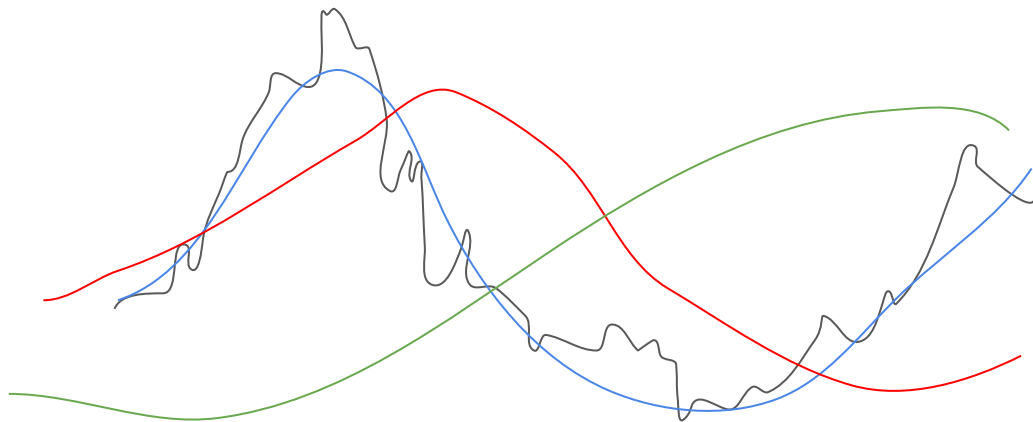
A heatmap showing the relevance of date features (year, month, day, day_of_week) to a target variable (sales). The 'year' feature has a high relevance of 0.978, highlighted in red. The 'month' feature has a relevance of 0.067, 'day' has 0.005, and 'day_of_week' has 0.006. The 'sales' column is highlighted with a blue border, and the values 0.081 (for year) and 0.037 (for day_of_week) are circled in yellow. A speech bubble points to the 'sales' column with the text 'More relevant than other date data.'

year -	0.978	0.000	0.000	0.081	0.
month -	0.067	-0.000	0.000	0.020	0.
day -	0.005	-0.000	0.000	-0.012	0.
day_of_week -	0.006	0.000	0.000	0.037	-0.
	id -	store_nbr -	family -	sales -	

More relevant than
other date data.

Enhanced Lag Features

- Used multiple lags: 1 day, 7 days, 30 days
- Applied Exponentially Weighted Moving Average (EWMA)
- Captured sales trends more accurately



Data Combination and Cleaning

- Combined training and test data for consistent feature engineering
- Cleaned data to remove noise

Conclusion

- Reduced features to simplify the model
- Prevented overfitting

The Models

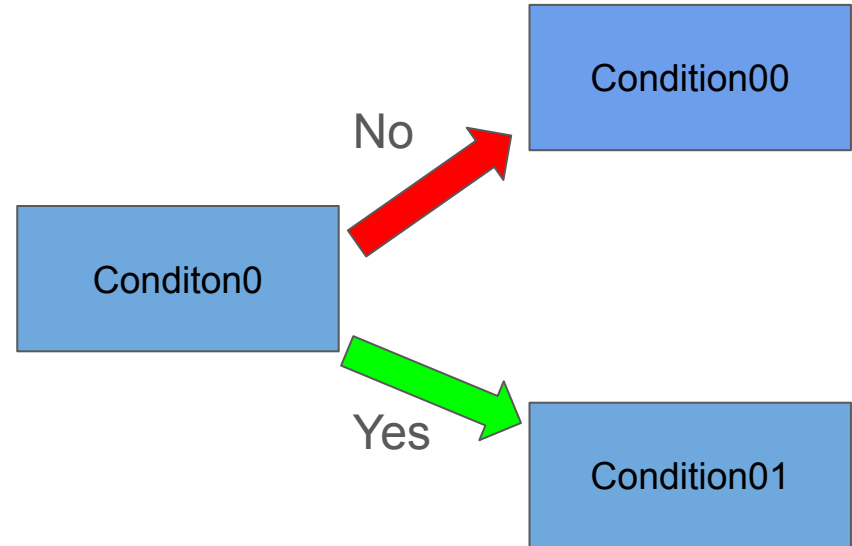
Things to keep in mind

- Handling discrete variables
- The target variable is continuous
- The data is a time series

The Basic : Decision Tree

Core characteristic :

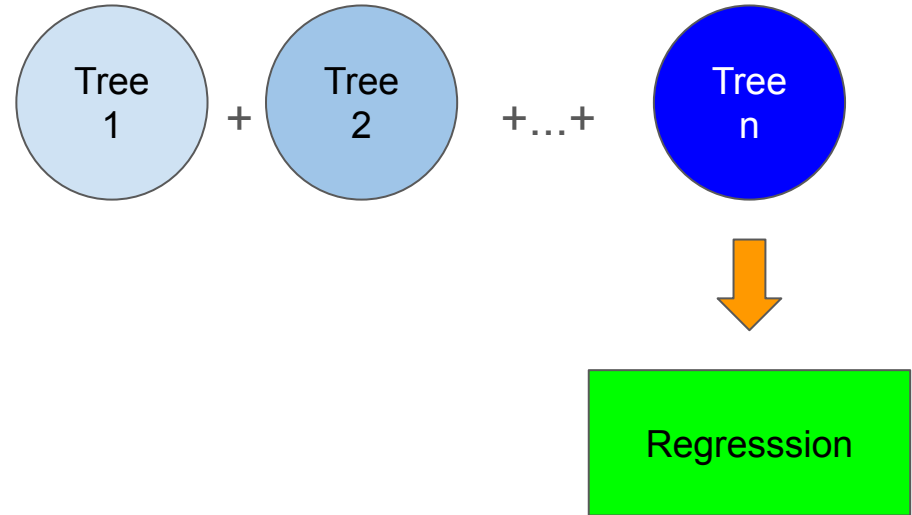
- **Good** at learning **discrete data**
- Extremely **prior to noises** in the data
- **Not robust** against a big amount of data
- Is **not** a **time series** model



Our Model: Light Gradient Boosting Machine Regressor

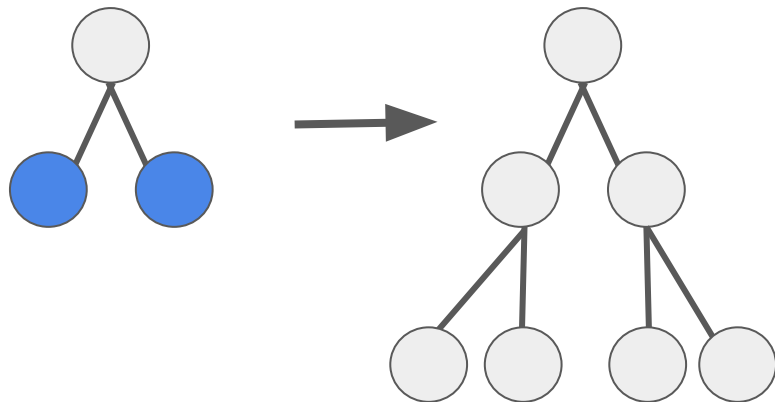
Core characteristic :

- **Good** at learning **discrete data**
- **Prior to noises** in the data
- Is **not** a **time series** model
- **Need hyperparameter tuning**

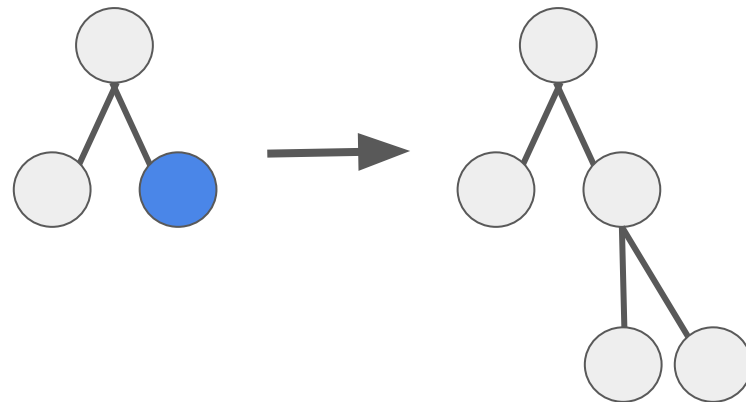


LGBM vs GBM

GBM



LGBM

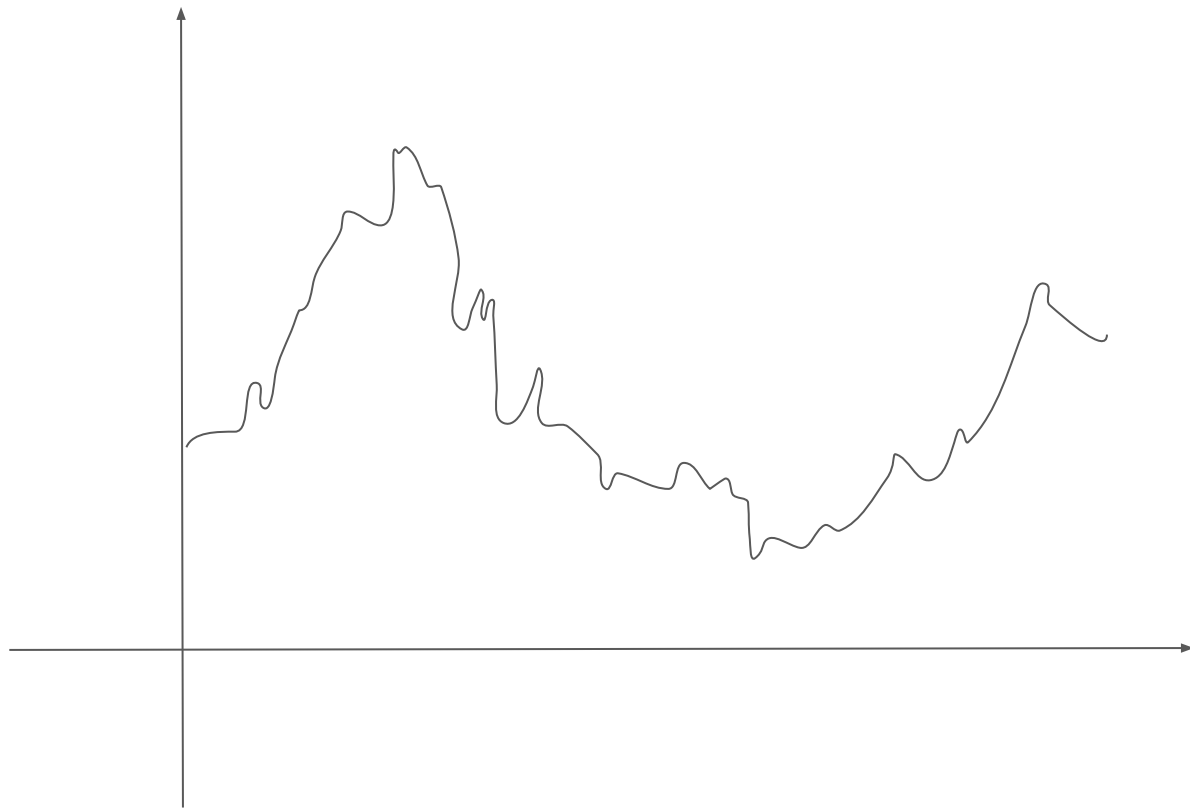


The Sliding Window

Sales:	Sales-1:	Sales-2:		Sales-n:
1	Nan	Nan		Nan
2	2	Nan		Nan
3	3	1	...	Nan
4	4	2		Nan
5	5	3		...
6	6	4		1

Before

The Sliding Window



Before

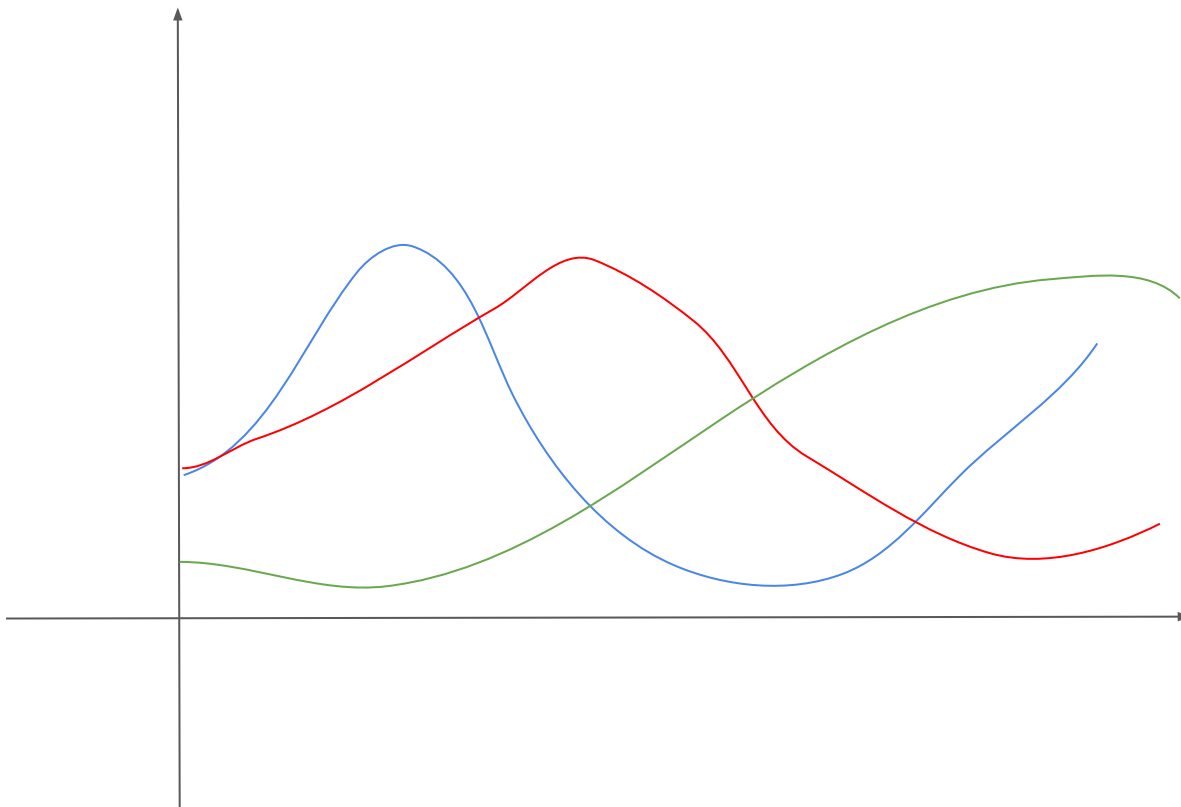
The Sliding Window

After

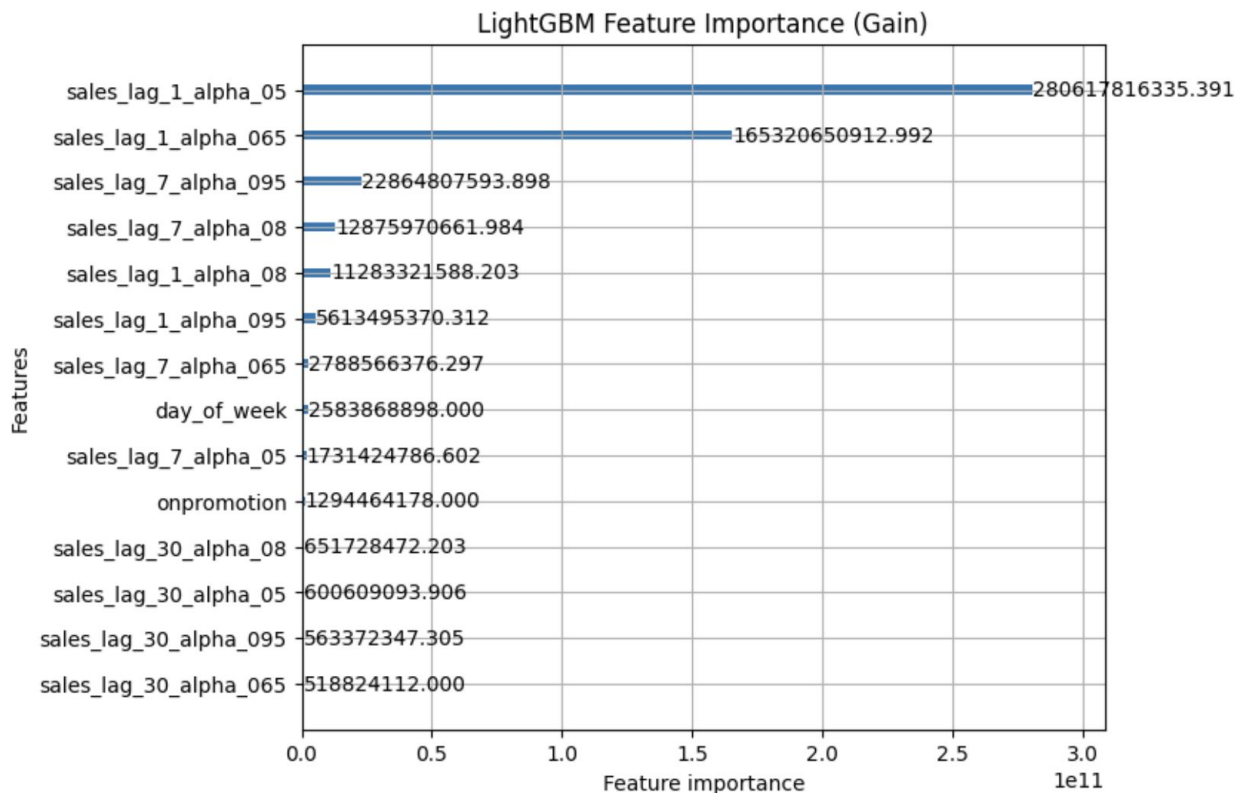
Sales:	Sales-1:	Sales-2:		Sales-n:
1	Nan	Nan		Nan
2	EMA(2)	Nan		Nan
3	EMA(3)	EMA(1)	...	Nan
4	EMA(4)	EMA(2)		Nan
5	EMA(5)	EMA(3)		...
6	EMA(6)	EMA(4)		EMA(1)

The Sliding Window

After



Gain of Each Features Ranked



Hyperparameter Tuning

Learning rate	: 0.1
Feature fracition	: 0.800087645
Bagging fraction	: 0.851134158
Bagging frequency	: 5
Verbose	: 0
Max depth	: 50
Num leaf	: 128
Max bin	: 512

What did we improve from before, and why?

- We **changed** our **model** from the regular **GBM** into **LGBM**.
- We **improve** our **sliding window**, which now gives better result
- We did **hyperparameter tuning**, the most important thing to do to **get better result** at GBM models.
- We **optimized** the **timeframe** used for training

Results

Our best result:

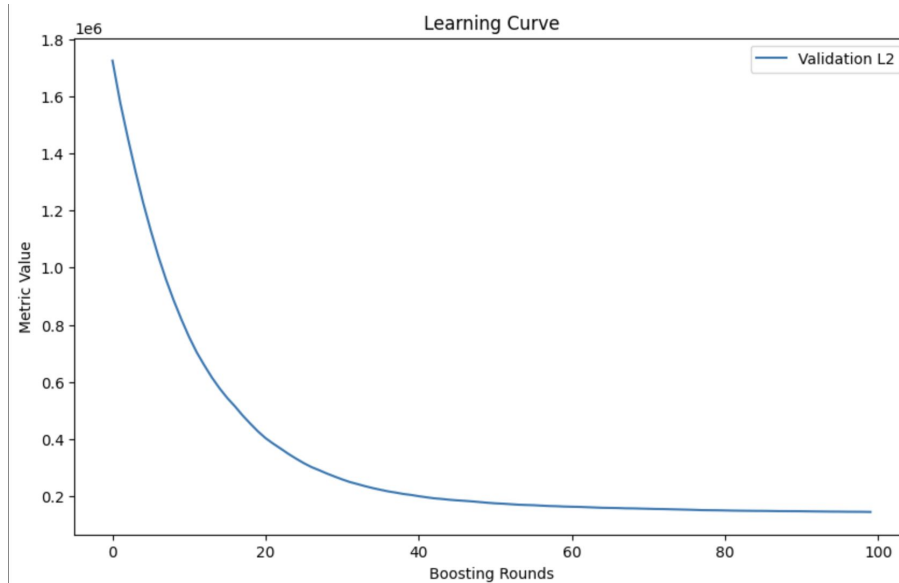
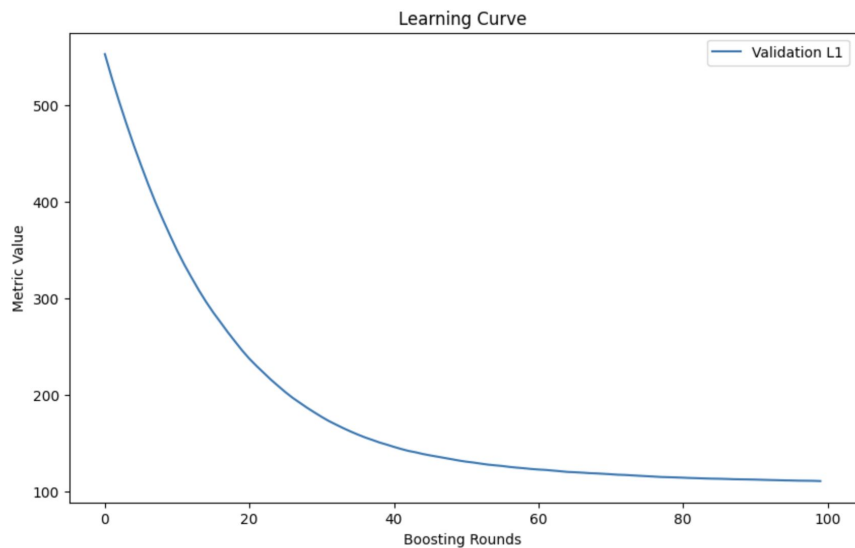


submission.csv

Complete · Ikhsan Rabbani · 11h ago · LGBM only train data without hp tuning

0.61289

Model Evaluation: Learning Curve



What did we learn?

Things We Notice....

- **Preprocessing** the data and exploratory data **analysis** are the **most important** thing to do before building a model
- **Sliding windows** can't just be added, but also needs to be **optimized** to give best prediction for time series data
- **Hyperparameter** needs to be **fine tuned** to give the best result

Thank You

Any Question?