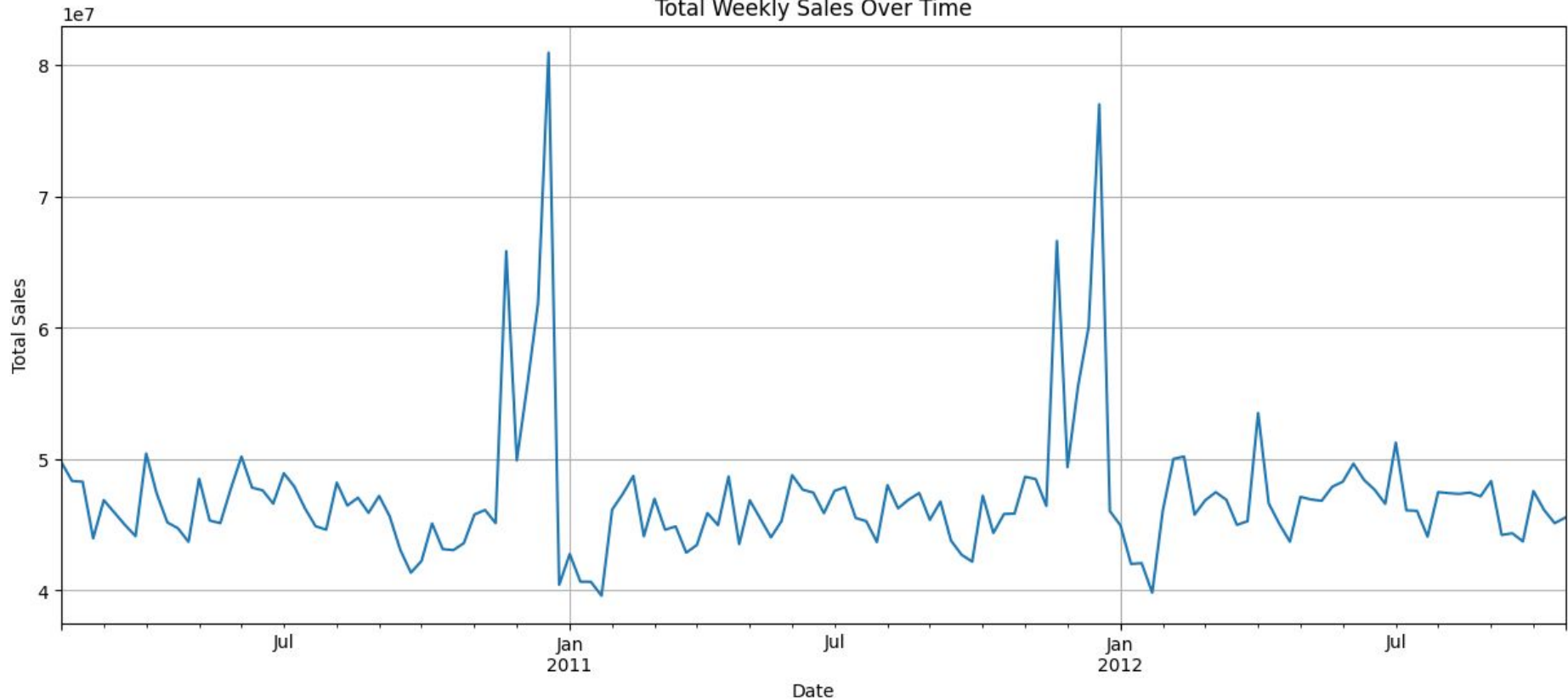


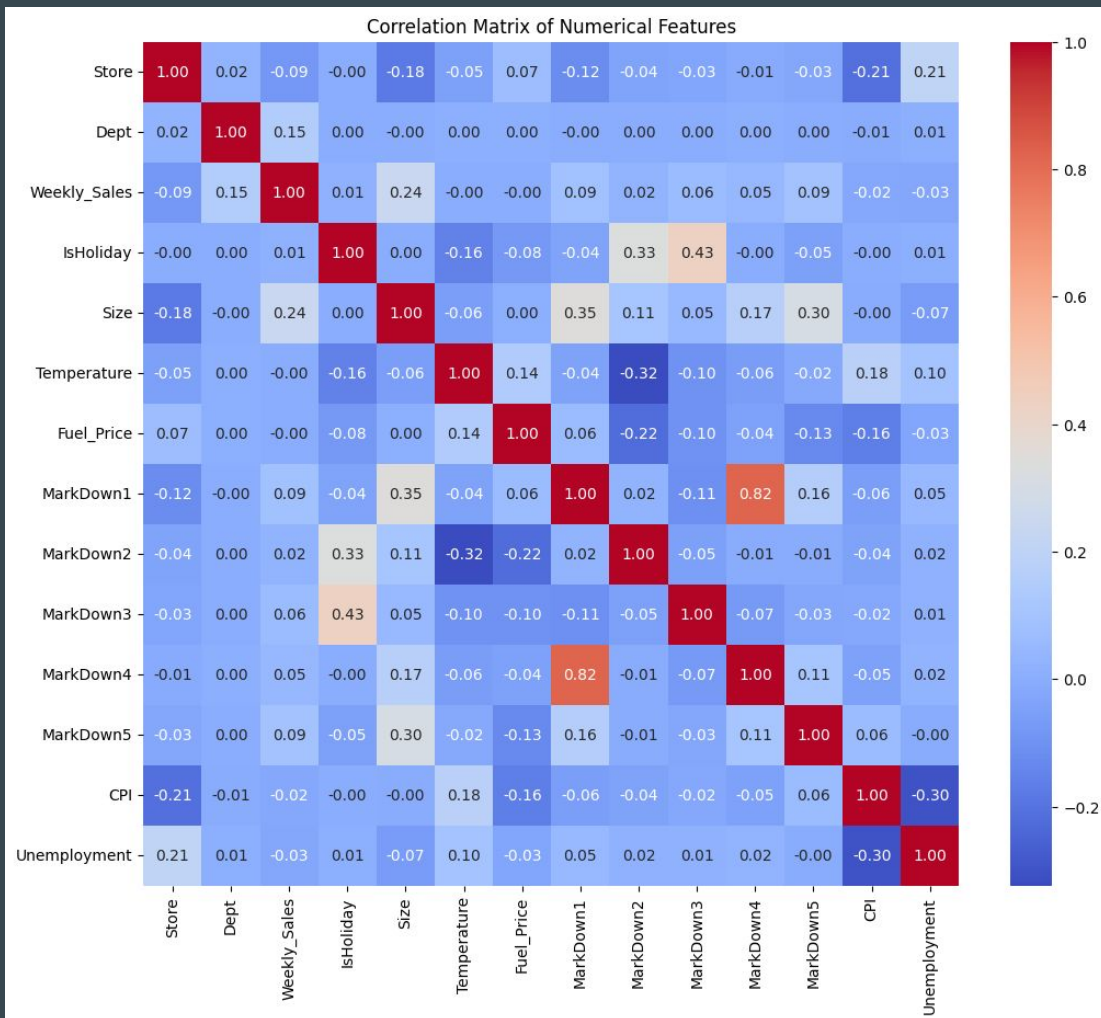
Total Weekly Sales Over Time



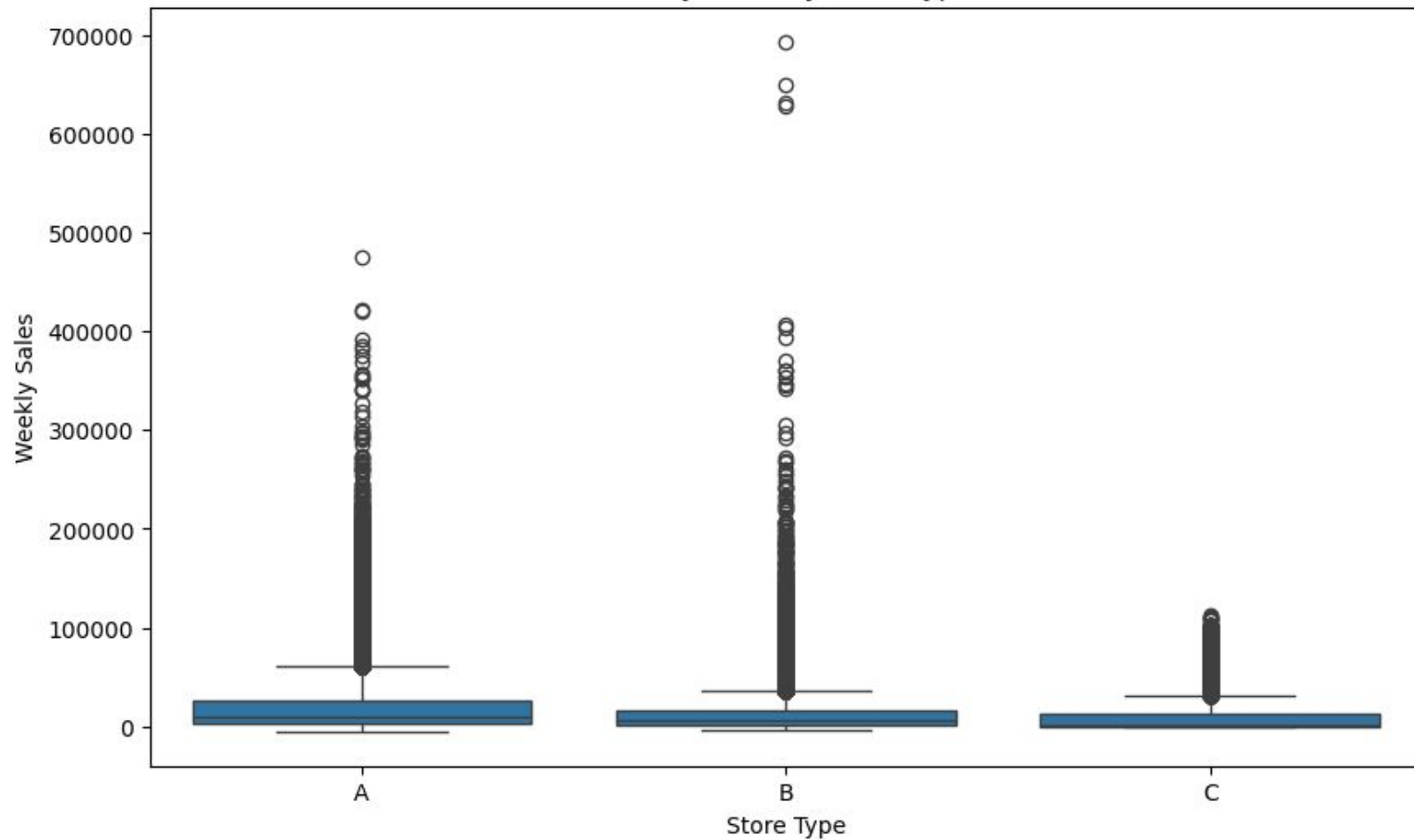
Walmart Recruiting - Store Sales Forecasting



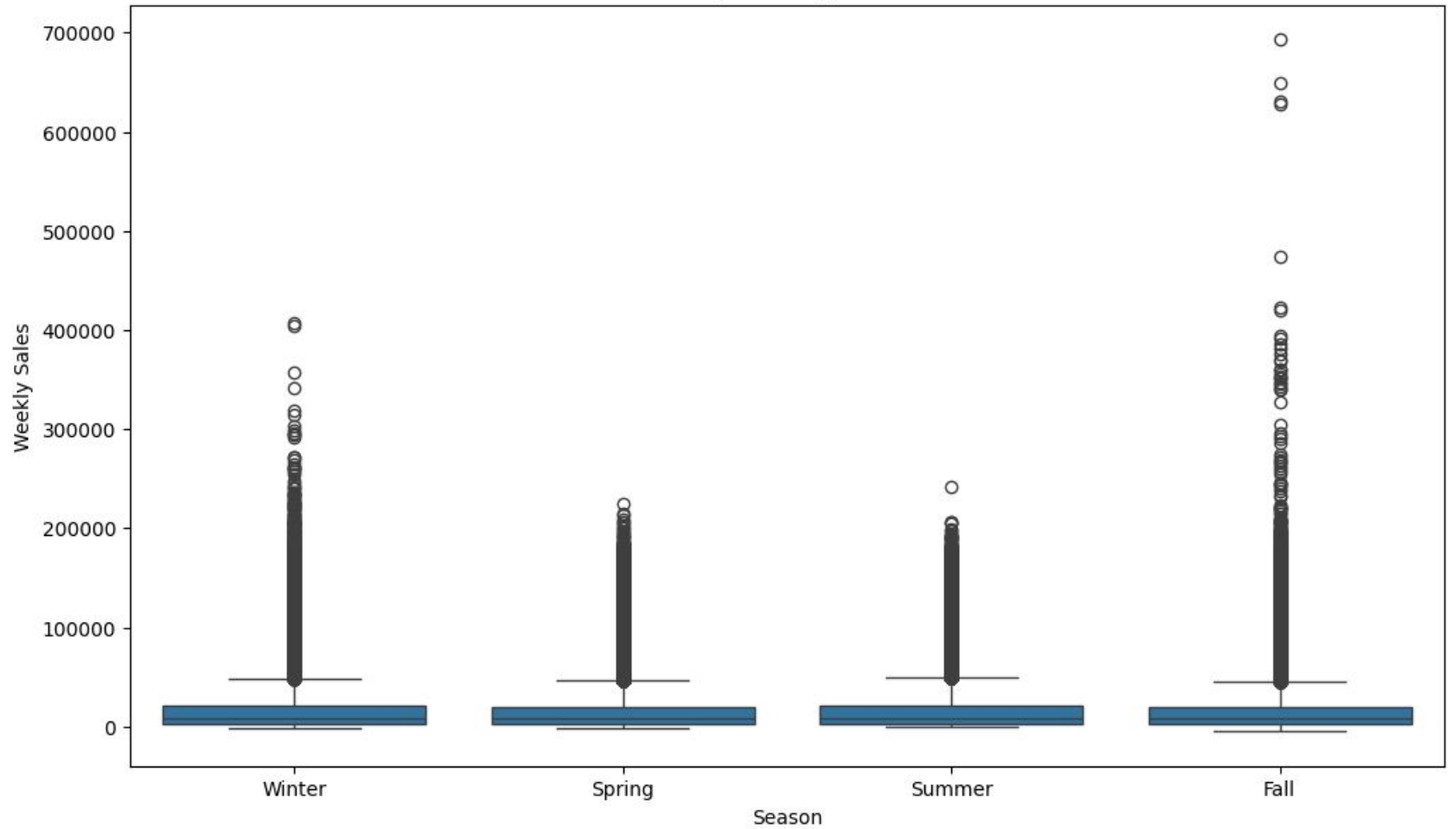
Mate Arevadze, Giorgi Toronjadze



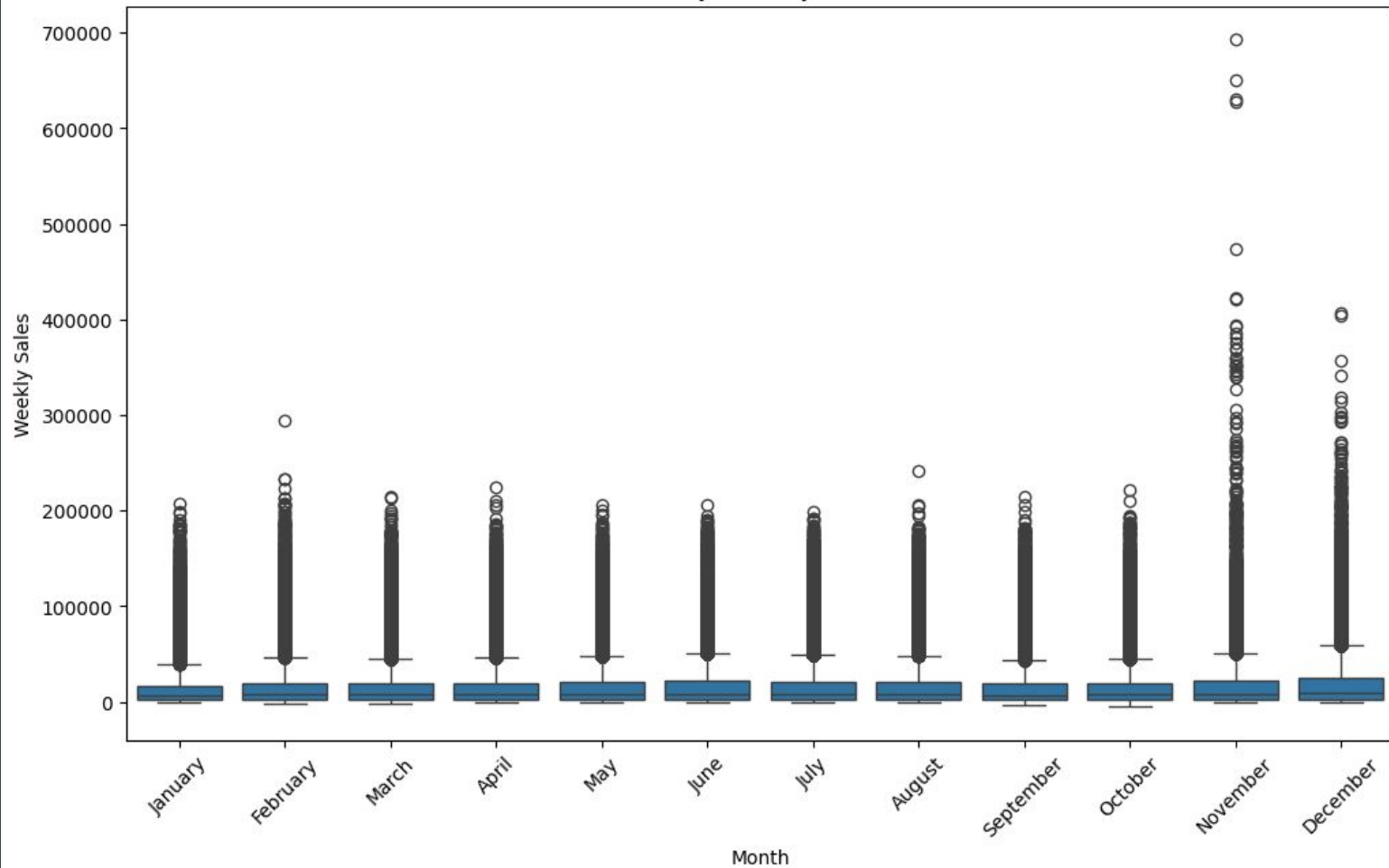
Weekly Sales by Store Type



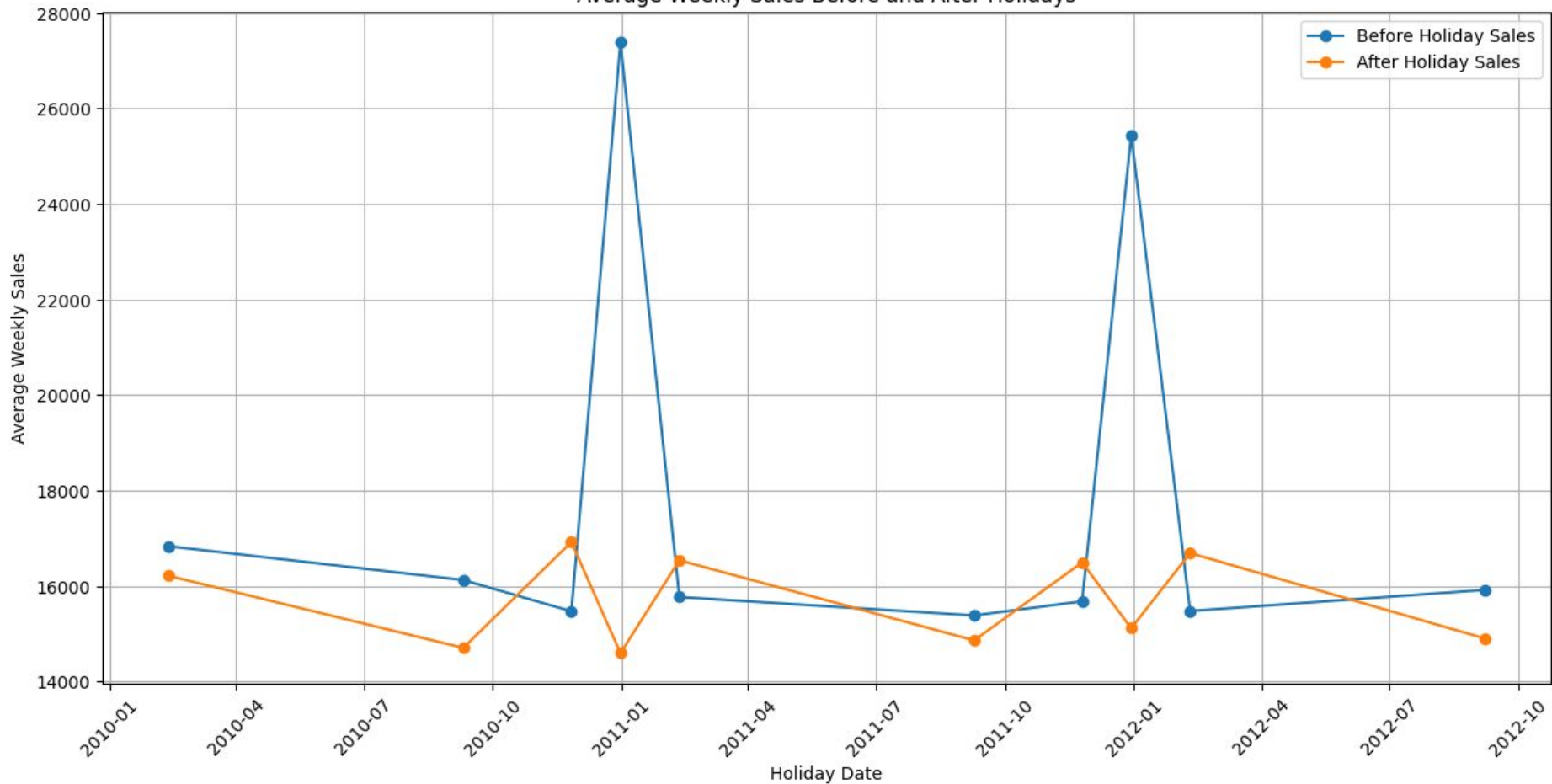
Weekly Sales by Season



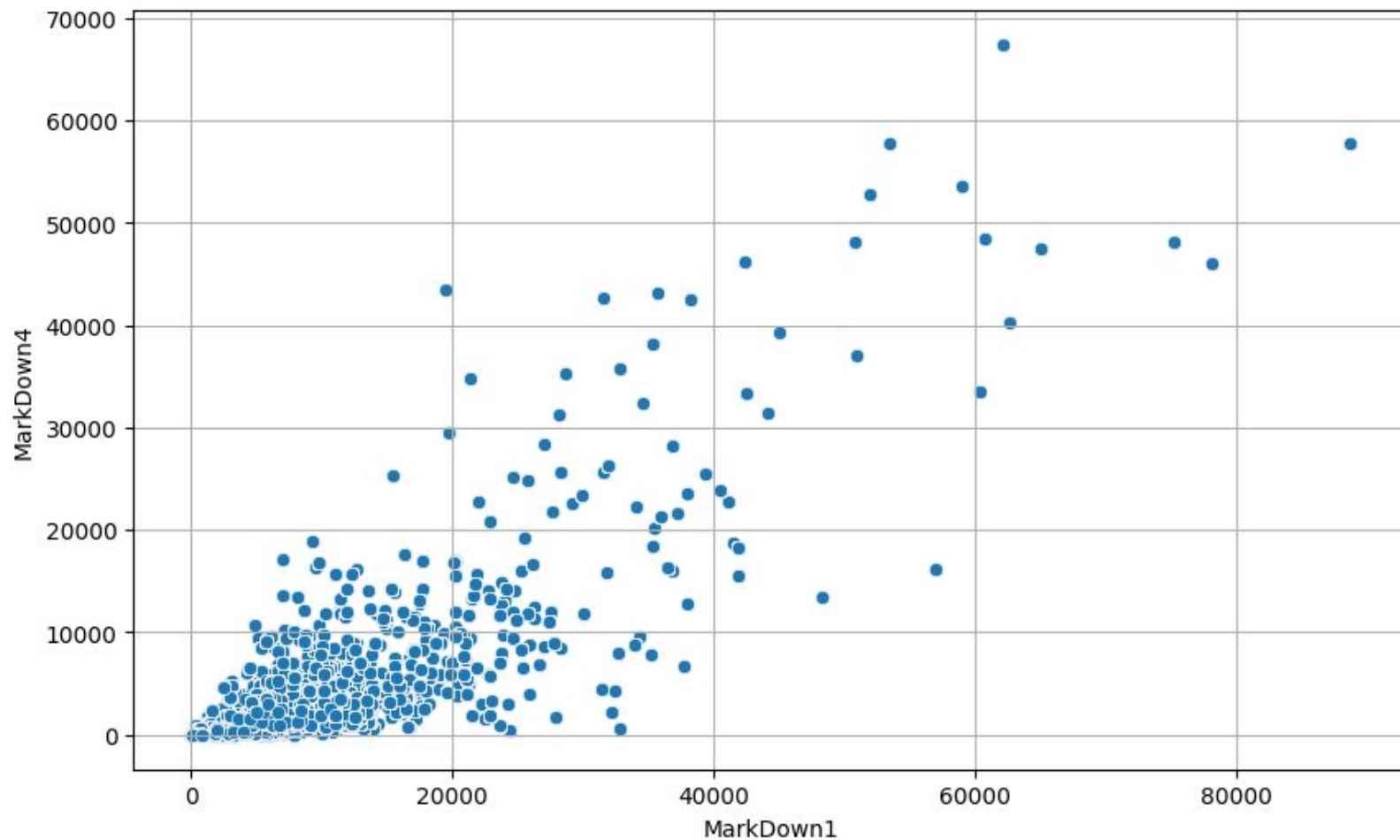
Weekly Sales by Month



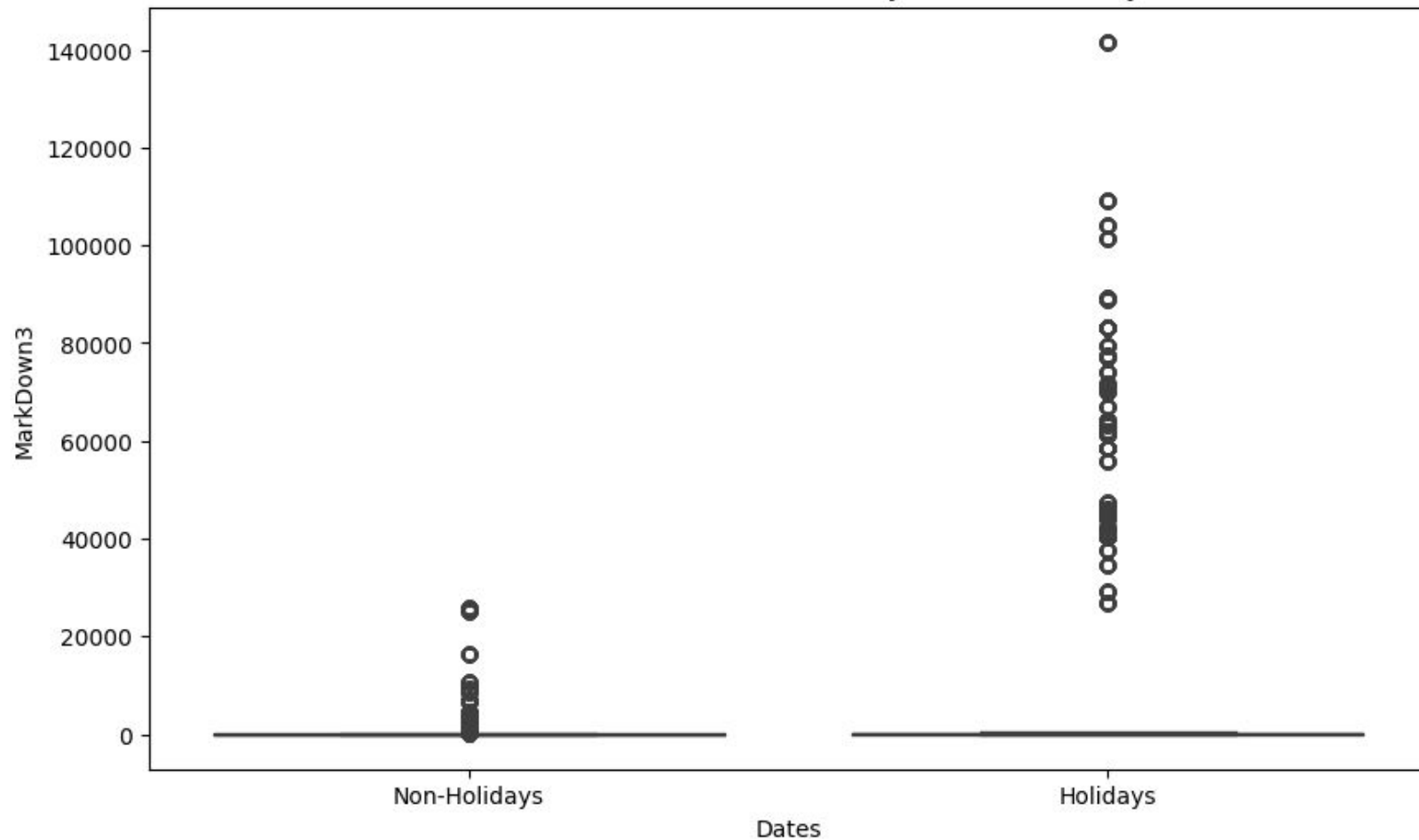
Average Weekly Sales Before and After Holidays



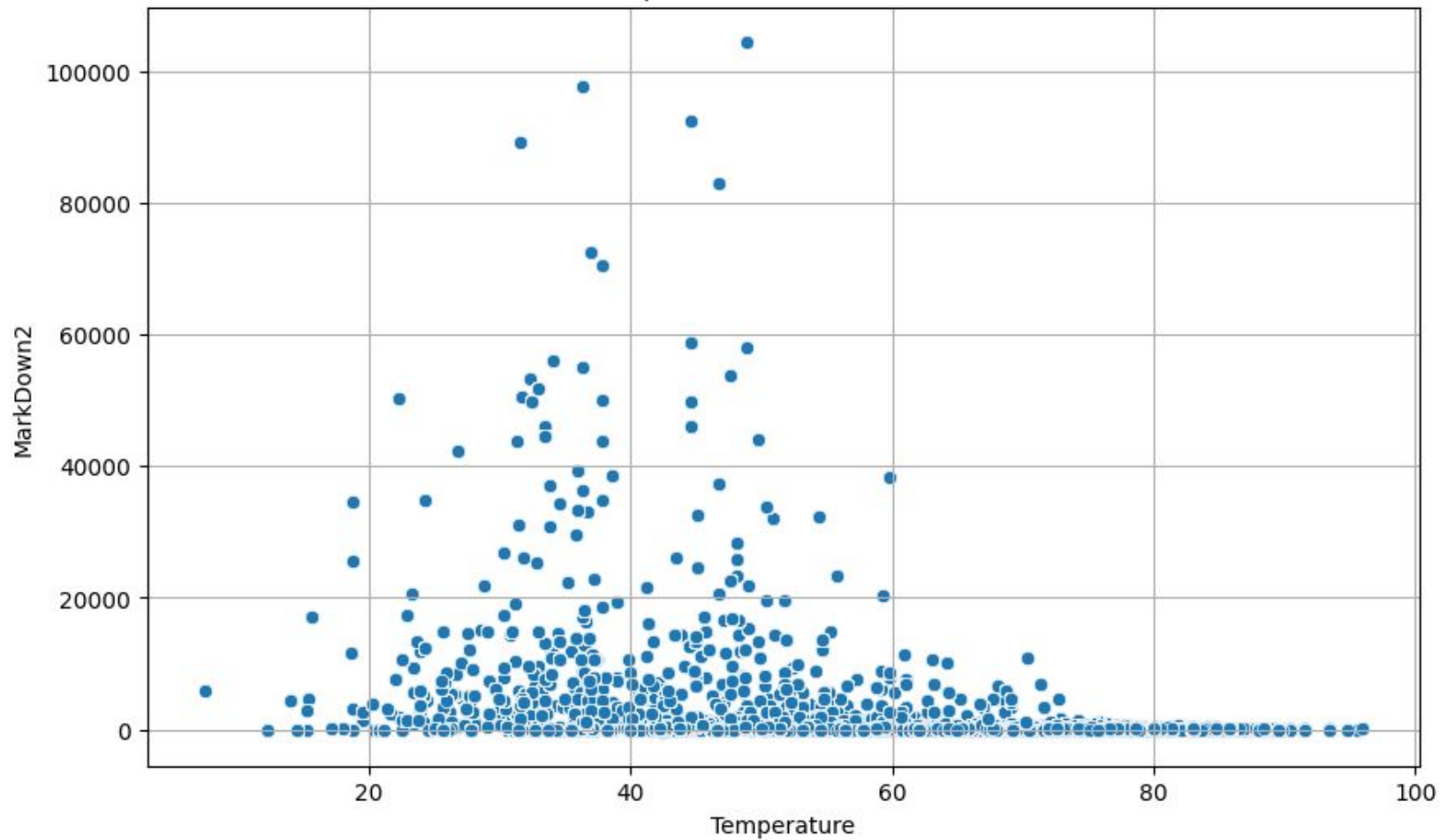
MarkDown1 vs. MarkDown4



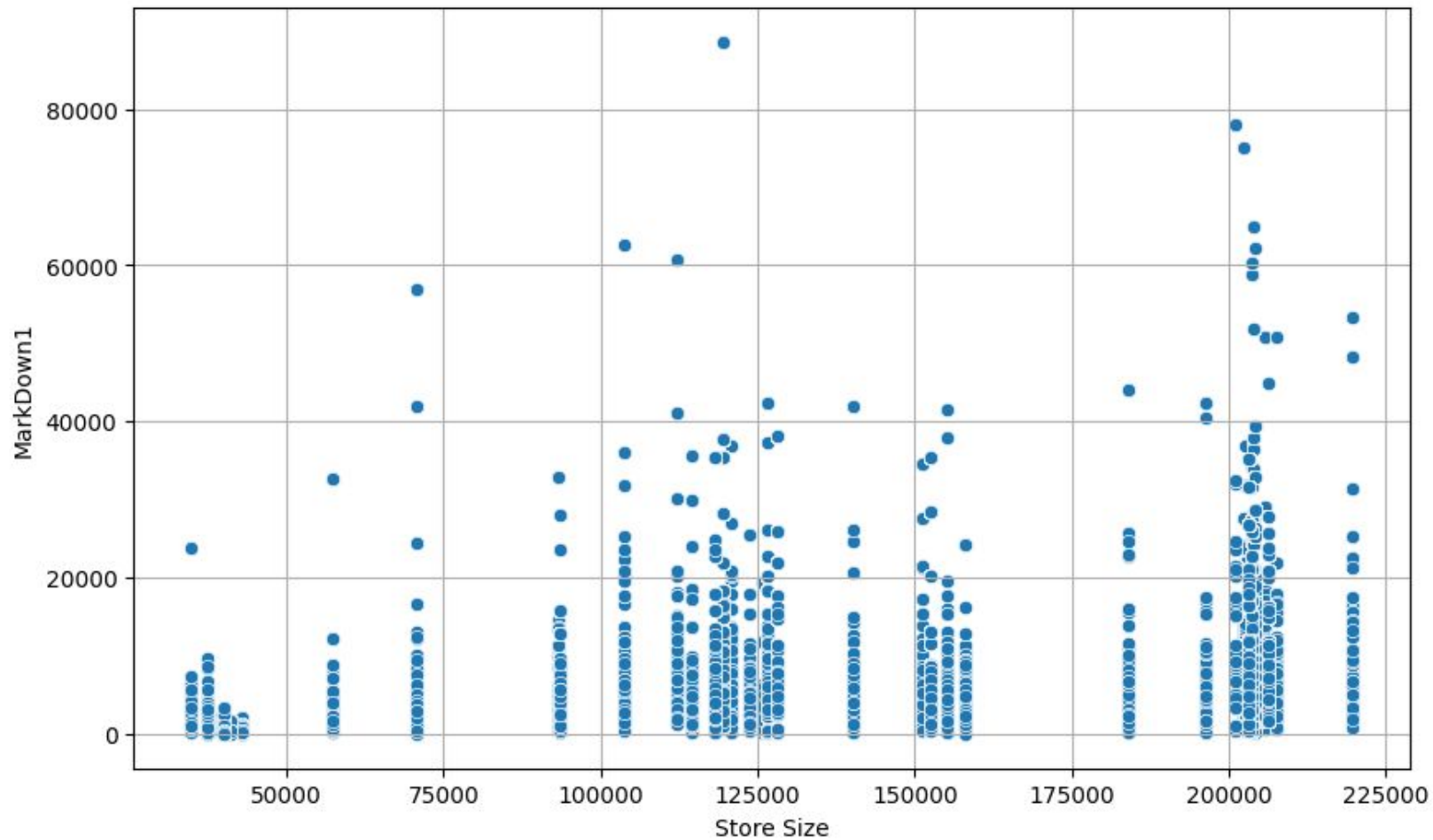
Markdown3 Distribution on Holidays vs. Non-Holidays



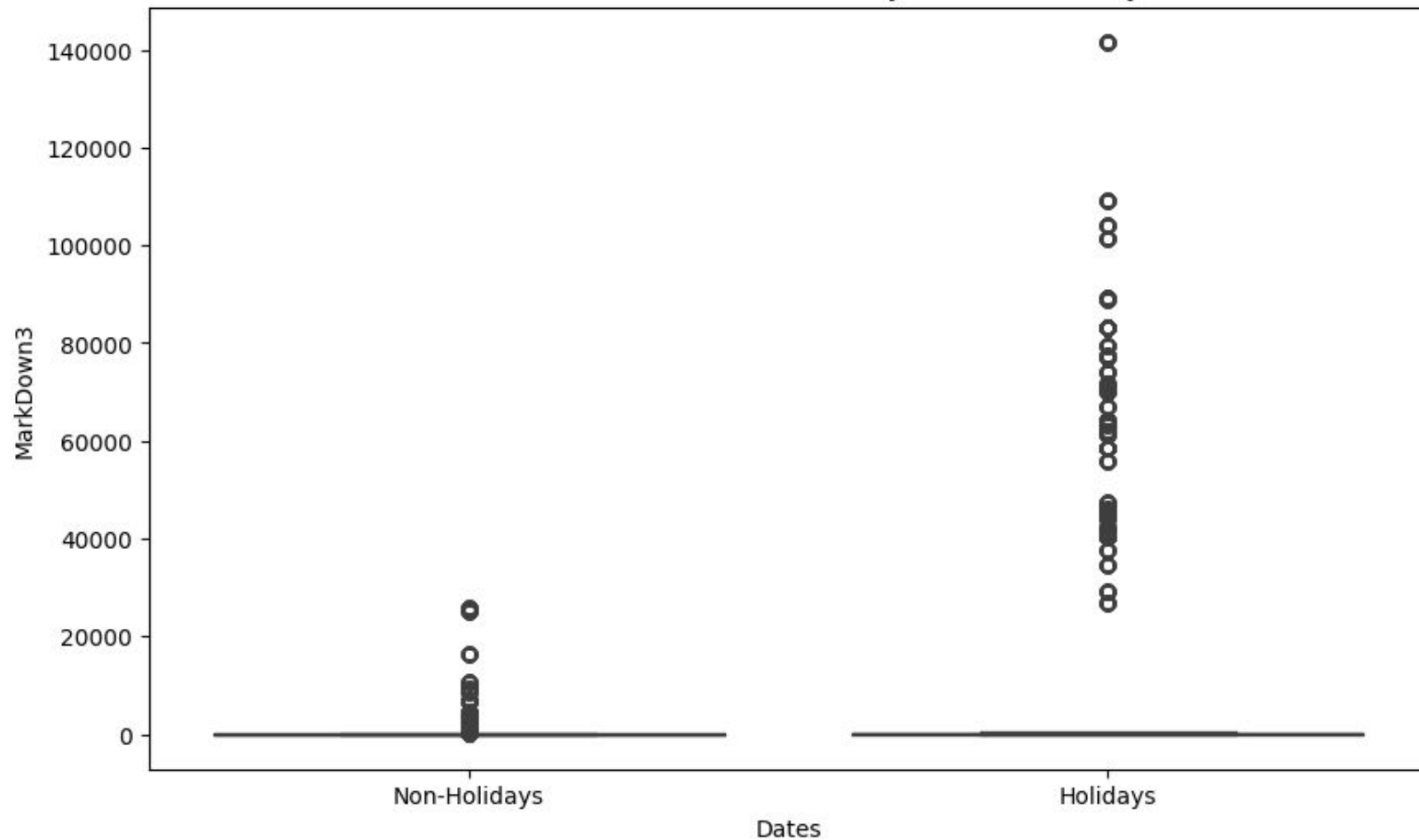
Temperature vs. MarkDown2



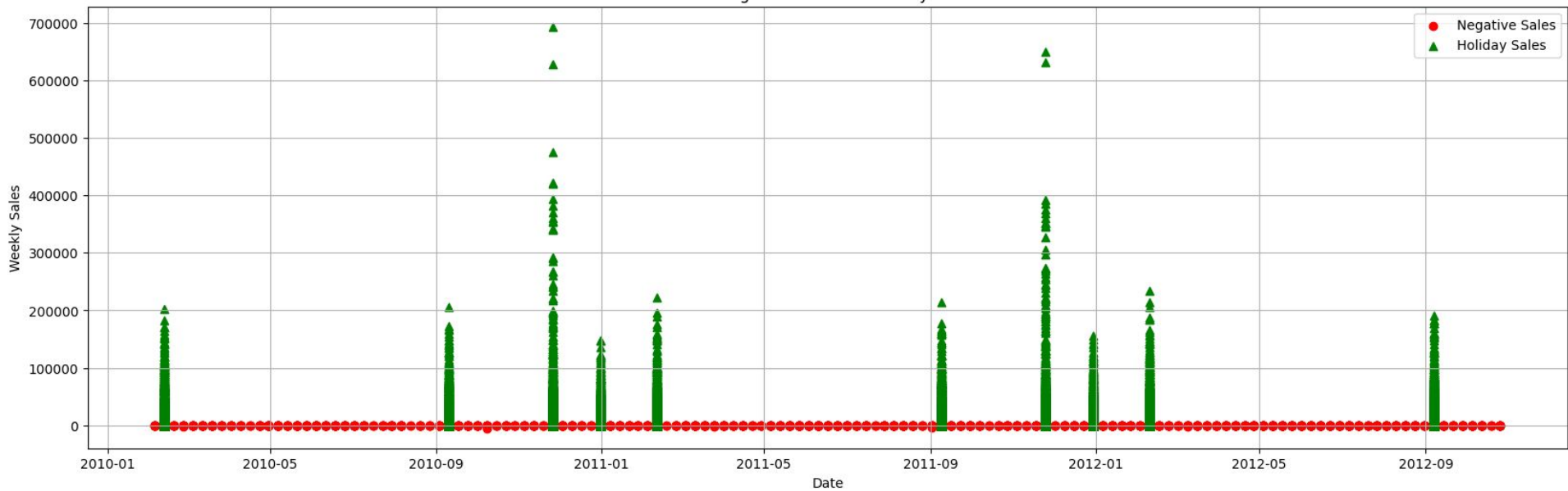
Store Size vs. MarkDown1



Markdown3 Distribution on Holidays vs. Non-Holidays



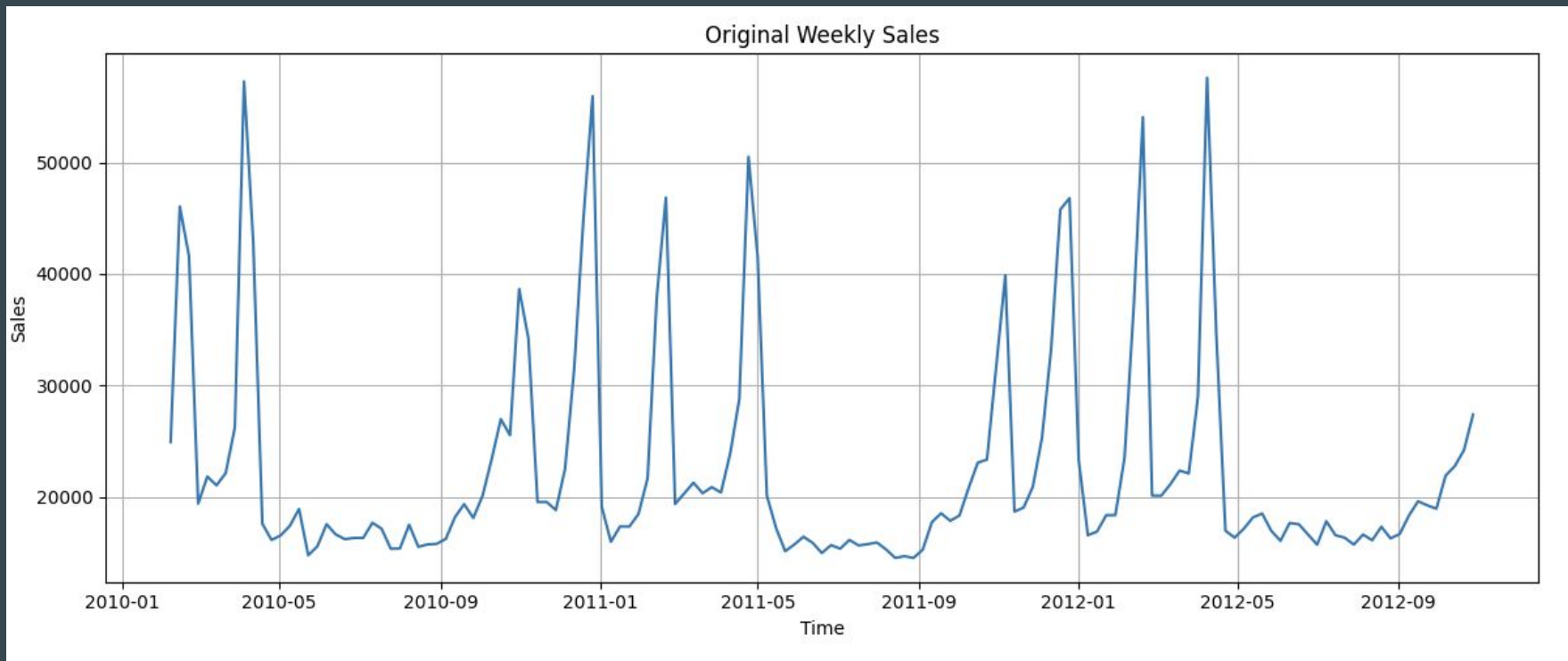
Negative Sales vs. Holidays



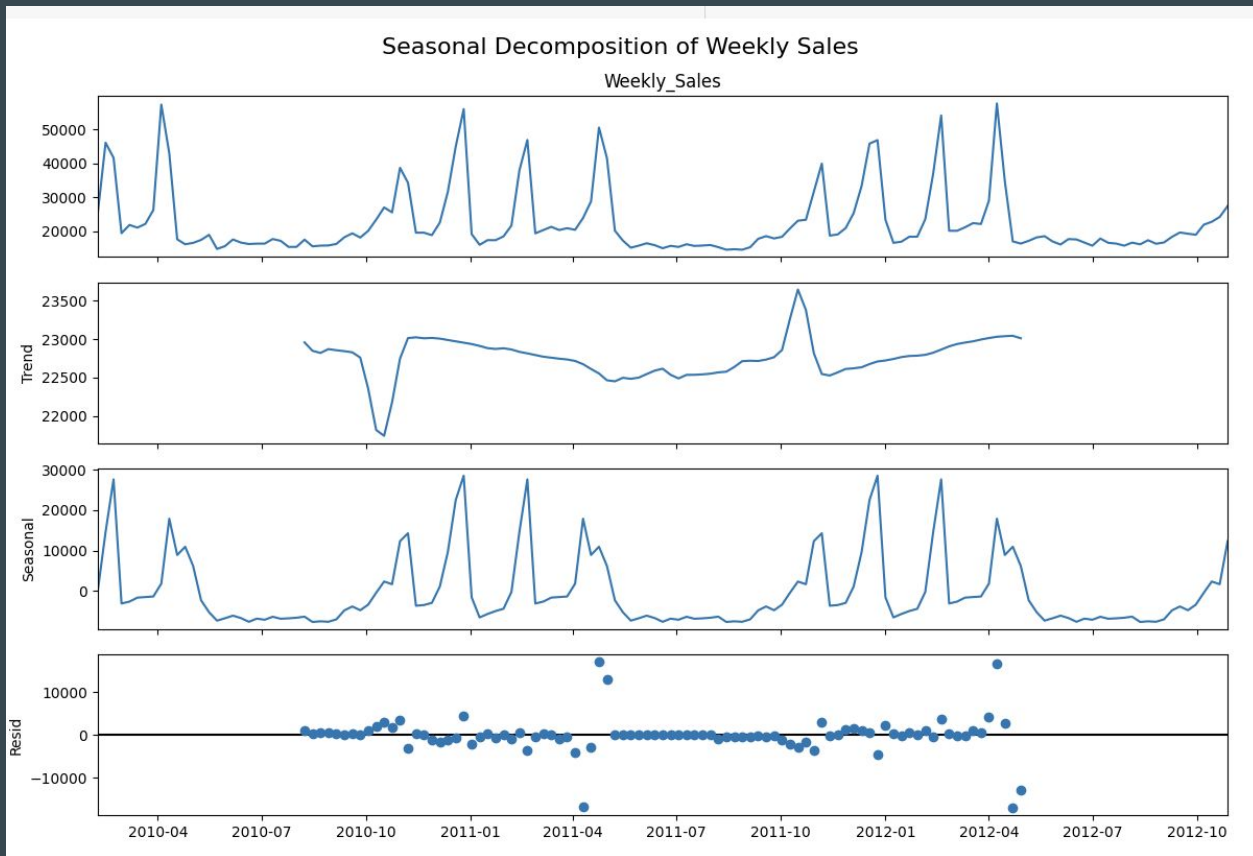
ARIMA/SARIMA/SARIMAX

- Used ARIMA (AutoRegressive Integrated Moving Average) for time series forecasting
- ARIMA(p,d,q): Models autoregressive (p), differencing (d), and moving average (q) components
- Challenges Encountered:
 - Non-stationarity: Sales data exhibited trends and non-constant variance.
 - Seasonality: Periodic patterns (e.g., holidays) not captured by basic ARIMA.
 - Errors: Inconsistent sample sizes, infinite values, and shape mismatches in some cases.

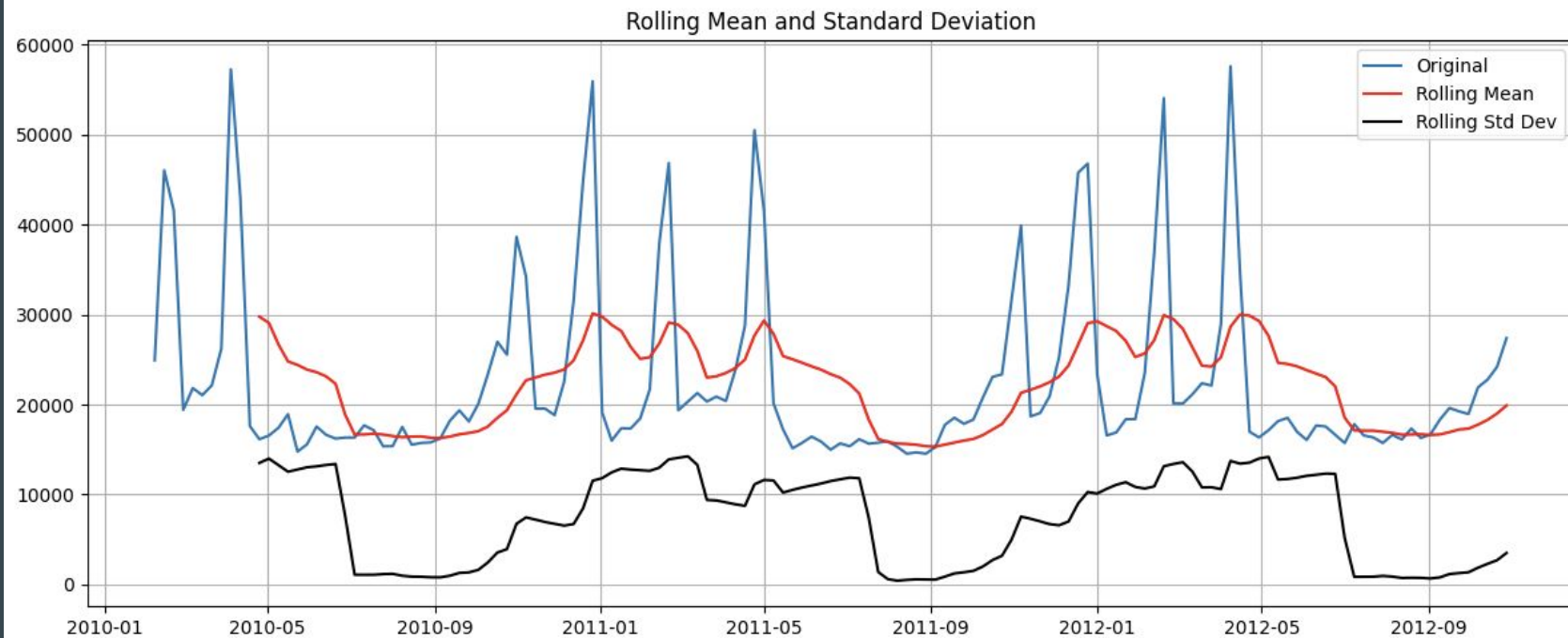
the sales data shows a clear upward/downward trend and fluctuating variance — meaning the data is non-stationary



Using seasonal decomposition, we can clearly observe repetitive seasonal patterns and a trend component



Both the rolling mean and standard deviation change over time



Exogenous Variables in SARIMAX Model

Exogenous variables (exogs) are external factors that impact sales.

Examples we use:

- **IsHoliday:** Whether the week includes a holiday (affects sales spikes).
- **Temperature:** Weather influences shopping behavior.
- **Fuel_Price:** Cost of fuel can affect store visits.
- **MarkDown1–5 :** Promotional discounts offered.
- **Season:** Season of the year.
- **Weeks_until_next_holiday:** Time until the next holiday.

Including these variables helps the model understand sales beyond historical patterns. This improves forecast accuracy.

PREPROCESSING

- Markdown columns had many missing values (64%-74% in some)
 - Fill with 0s missing Markdown Columns
 - Fill with mean other missing columns
-
- Got better result when filling with 0s

Hyperparameter Tuning with Grid Search

SARIMAX has several parameters:

- Non-seasonal: **p, d, q**
- Seasonal: **P, D, Q, s** (with seasonal period $s=52$ weeks)

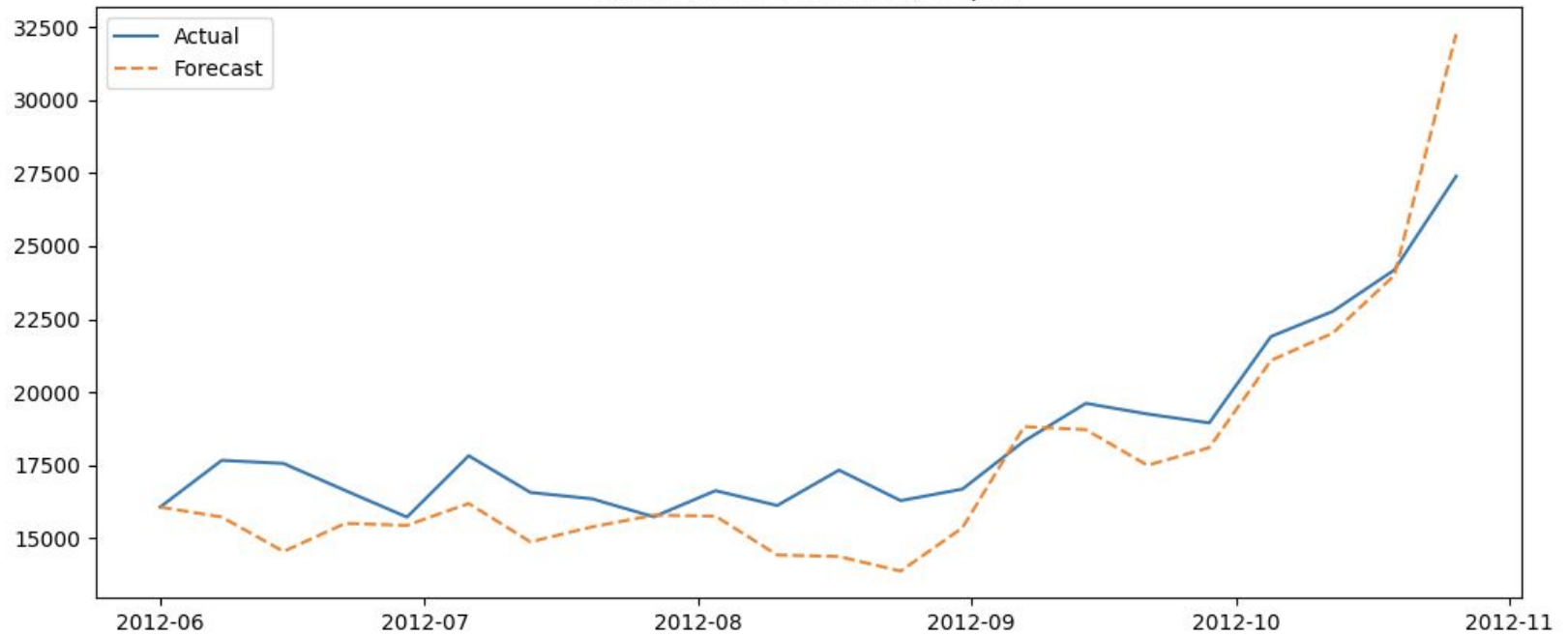
Choosing the right parameters is critical for model accuracy.

We automate this with **Grid Search**: trying all combinations in a defined range.

```
Best SARIMA: (1, 0, 1) x (0, 1, 0, 52) - AIC: 13.683230948927266
```

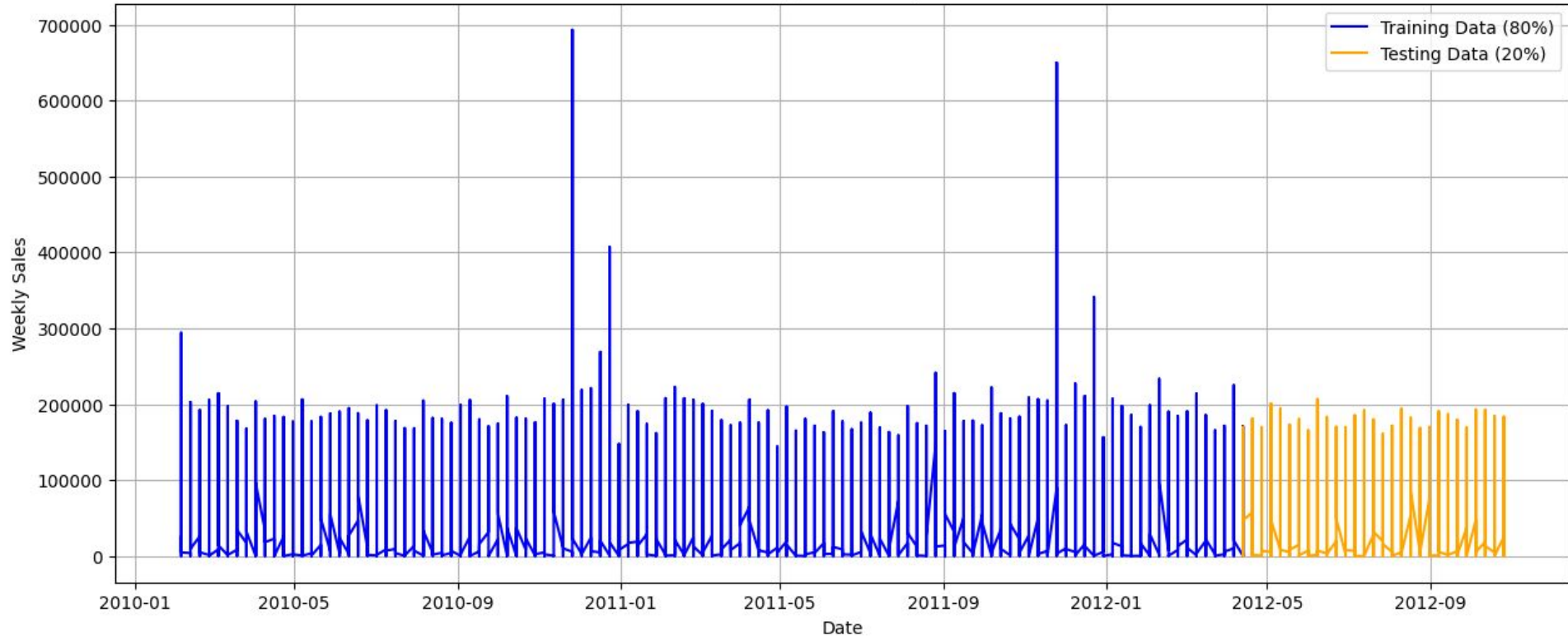
```
Best WMAE: 844.76
```

SARIMA Forecast: Store 1, Dept 1

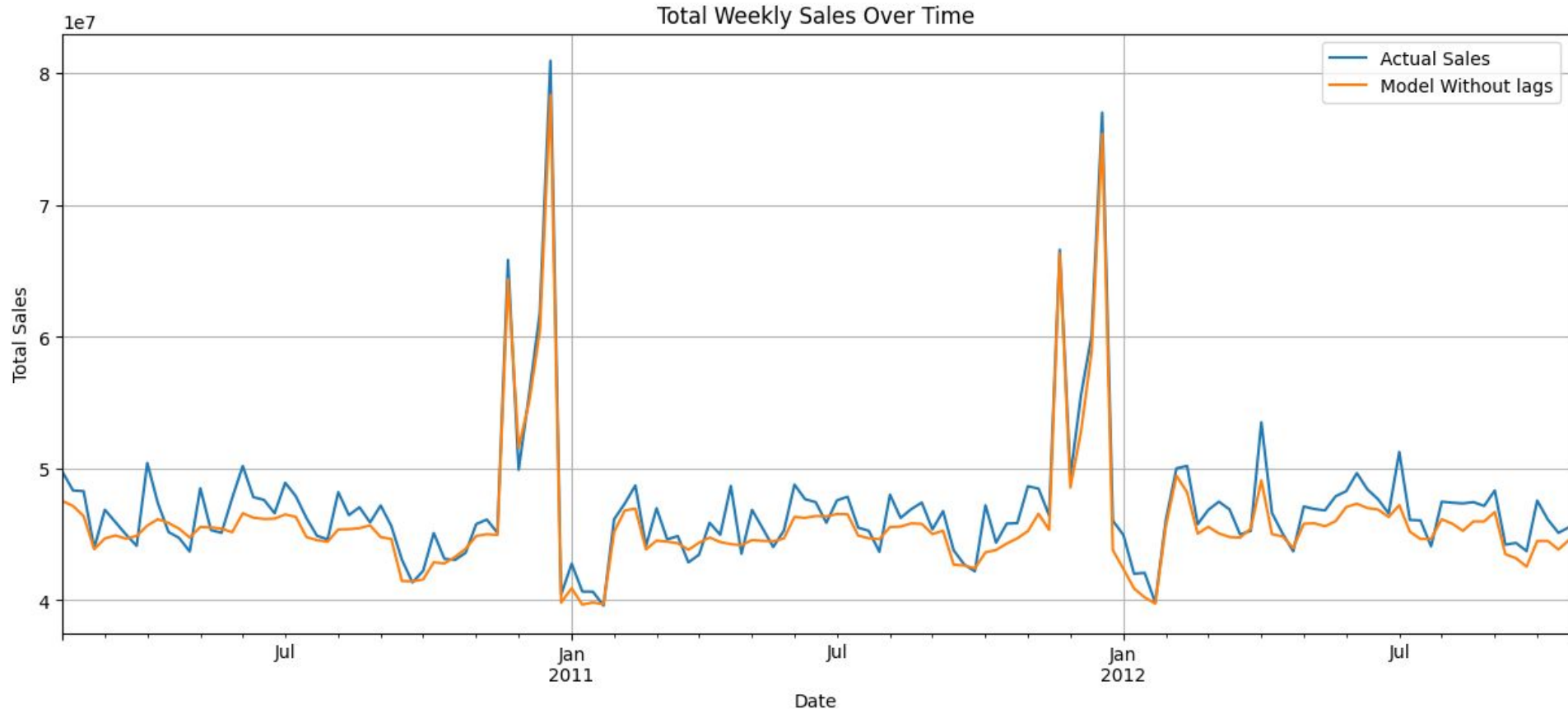


XGBOOST

Train/Test Data Split



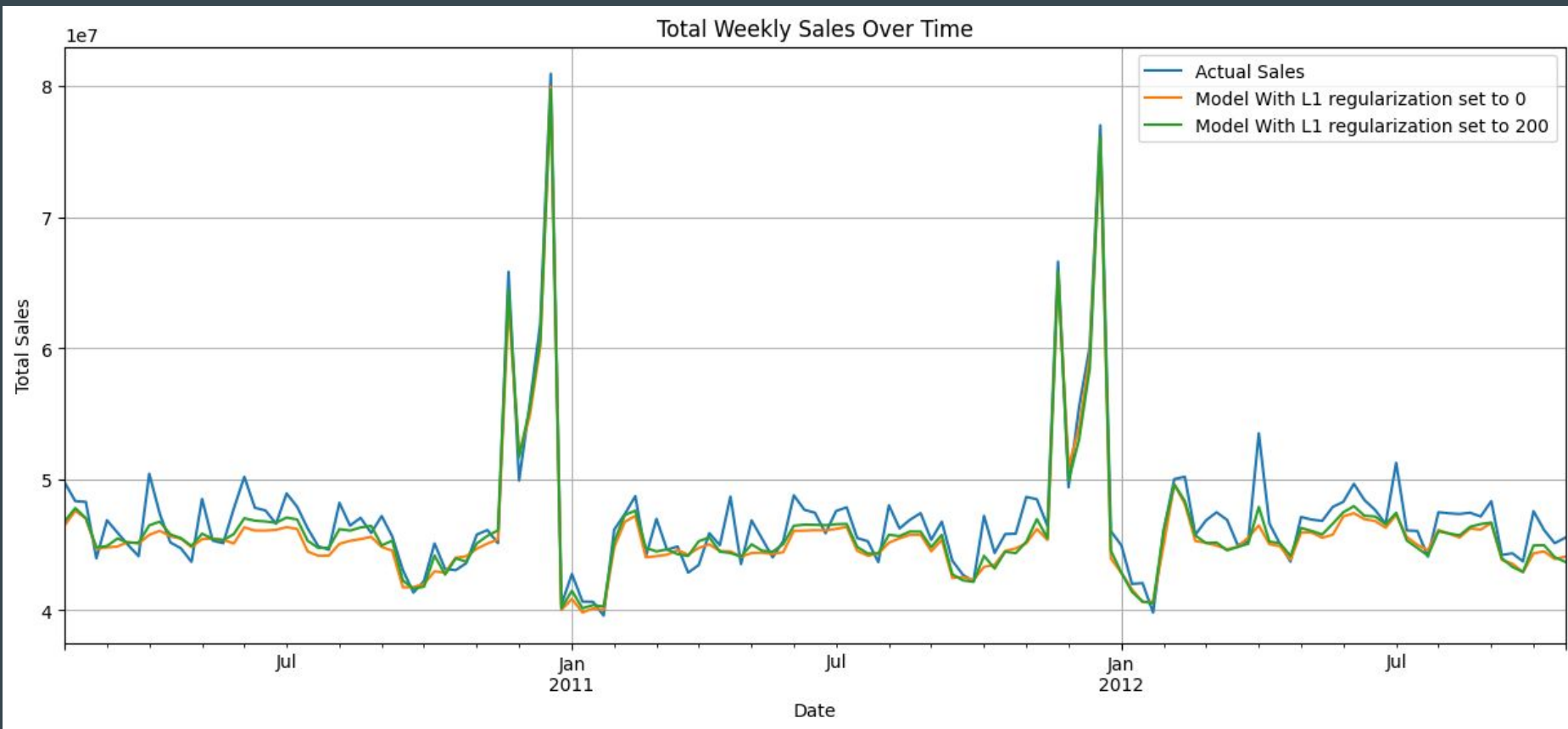
Without Lagging Features, Validation WMAE: 6436, With Lagging Features: 6000



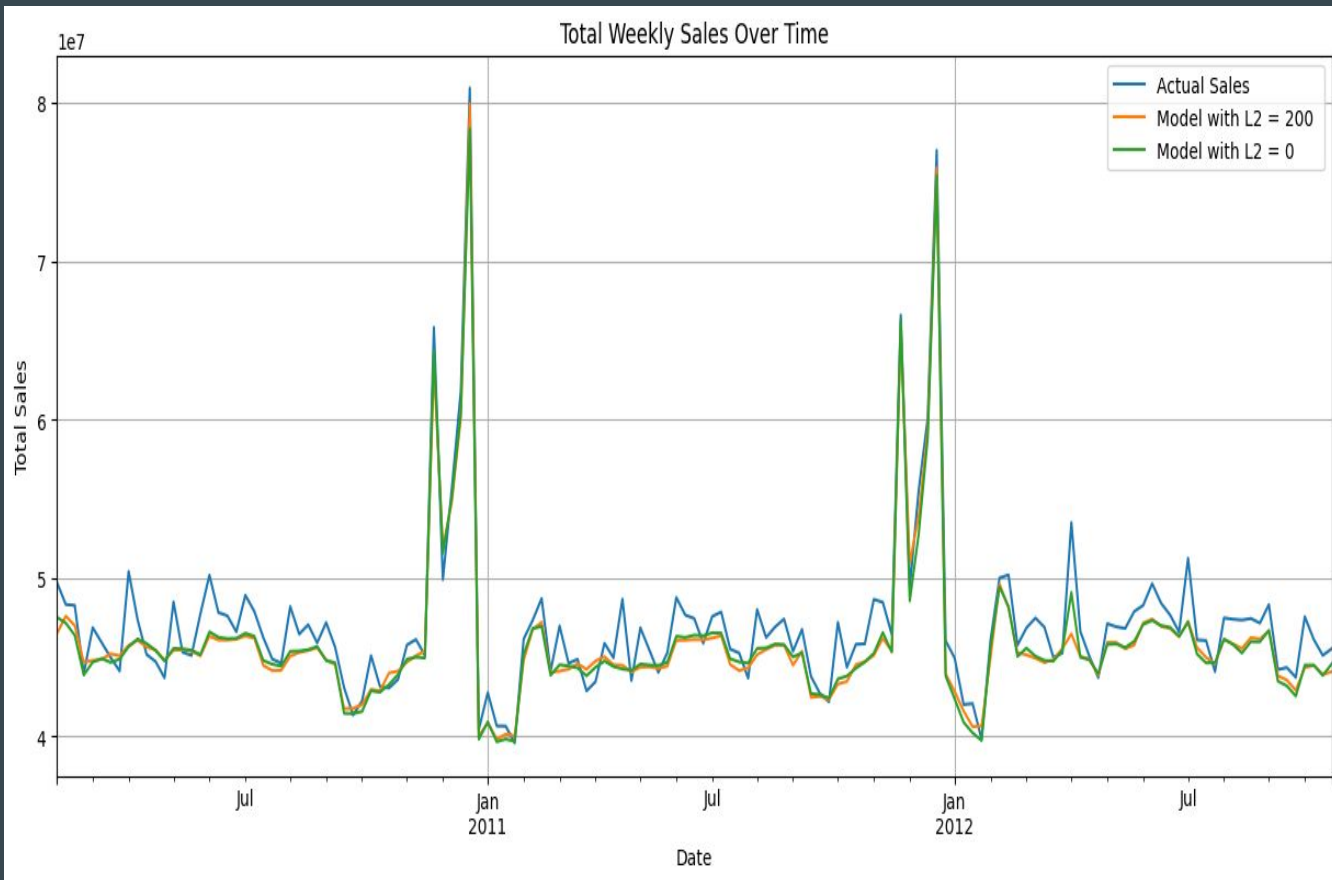
ALPHA REGULARIZATION EFFECTIVENESS

Validation WMAE for model with no L1 : 5700

Validation WMAE for model with L1 set to 200: 6400

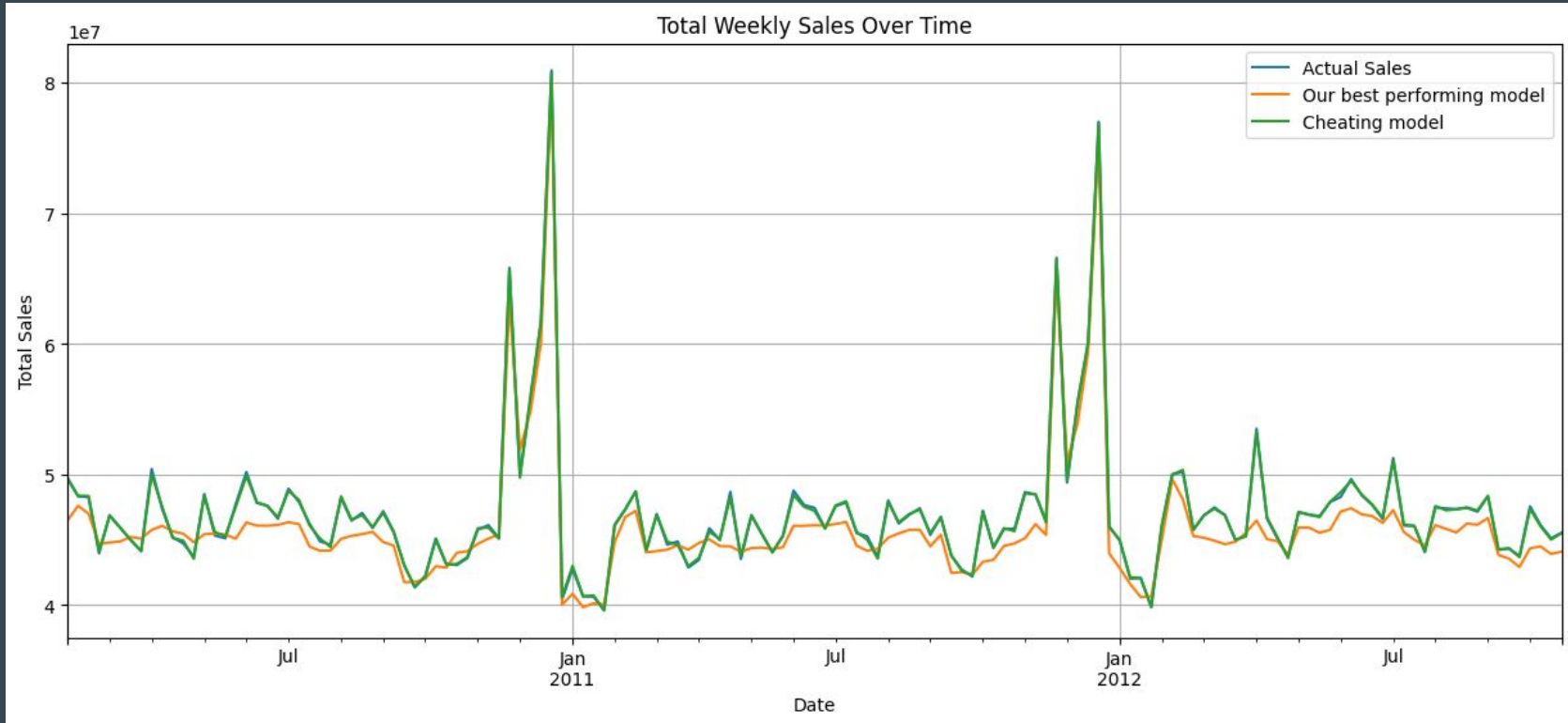


L2 (Lambda) Regularization effects



- ❖ Max_depth : 5
- ❖ n_estimators: 100
- ❖ reg_alpha: 0
- ❖ L2 Reg = 0: 6450 WMAE
- ❖ L2 Reg = 200 : 5800 WMAE

We Tried to see behavior of a “cheating” model that used random train/test split

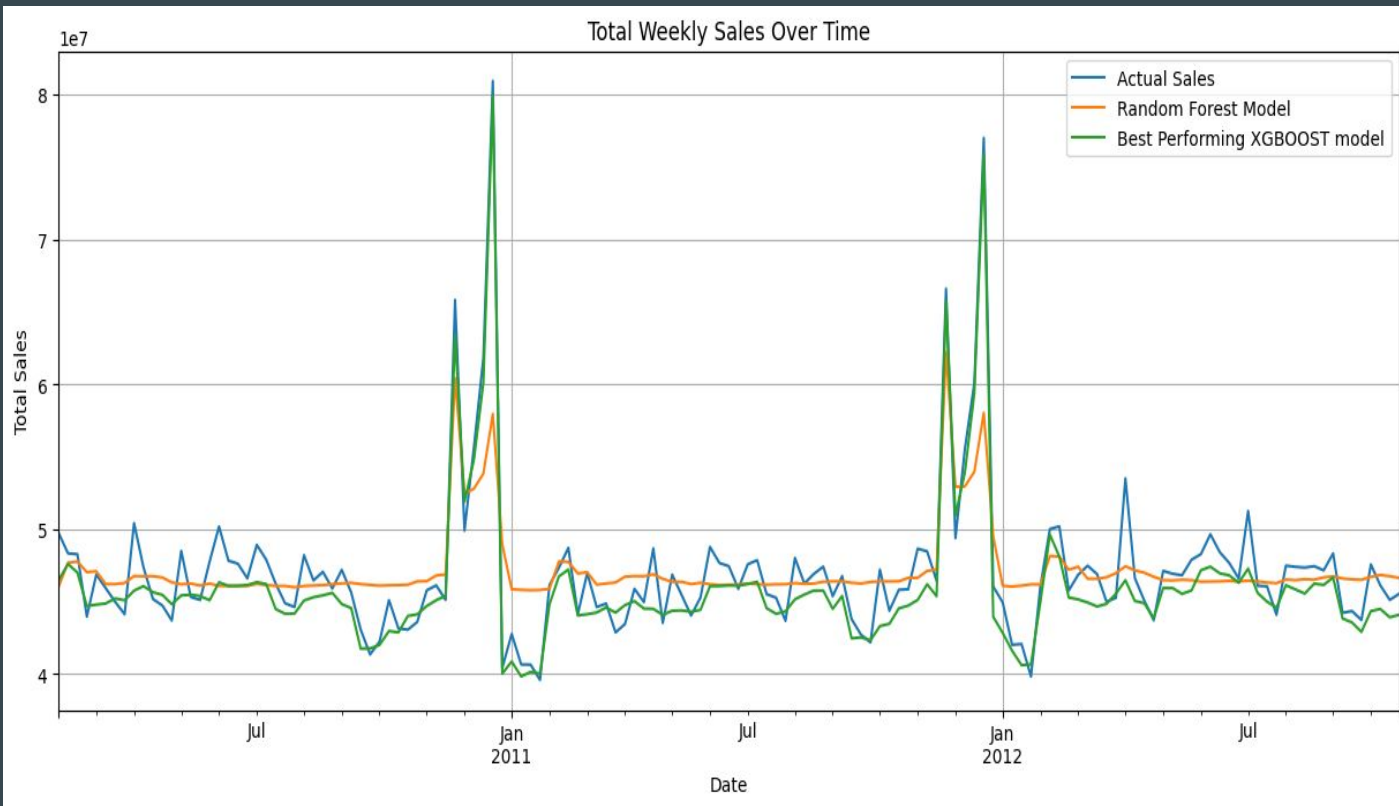


Other Features

- Max Depth was also an another important feature. setting it to high amounts (10,12) would cause overfitting, it would give us excellent results on train data (approx. 1800-2000), but it had bad results with validation data. even lambda regularization wouldn't work well if we didn't set max depth on conservative levels (5,6,7)
- Number of Estimators didn't affect the model that much, higher amount would improve the model, but not by significant margin
- we didn't experiment with learning rate that much, most of the time it was set to 0.1 or 0.3

RANDOM FOREST

Best Performing Random Forest Model compared to best performing XGBOOST model



Random Forest Model
WMAE: 5285

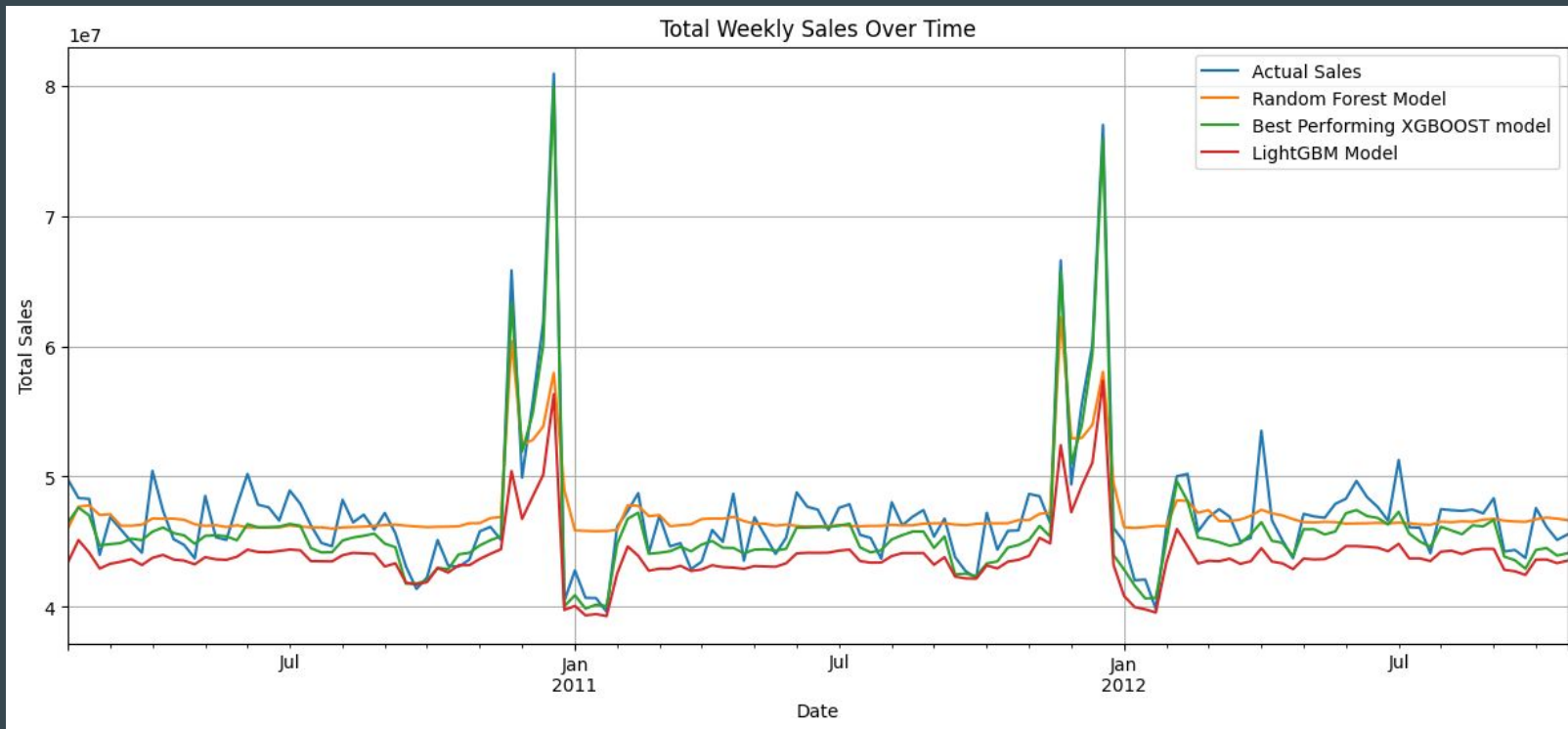
XGBOOST Model
WMAE: 5863

Despite this, XgBoost had a better score on kaggle and it was better at detecting seasonality and trends

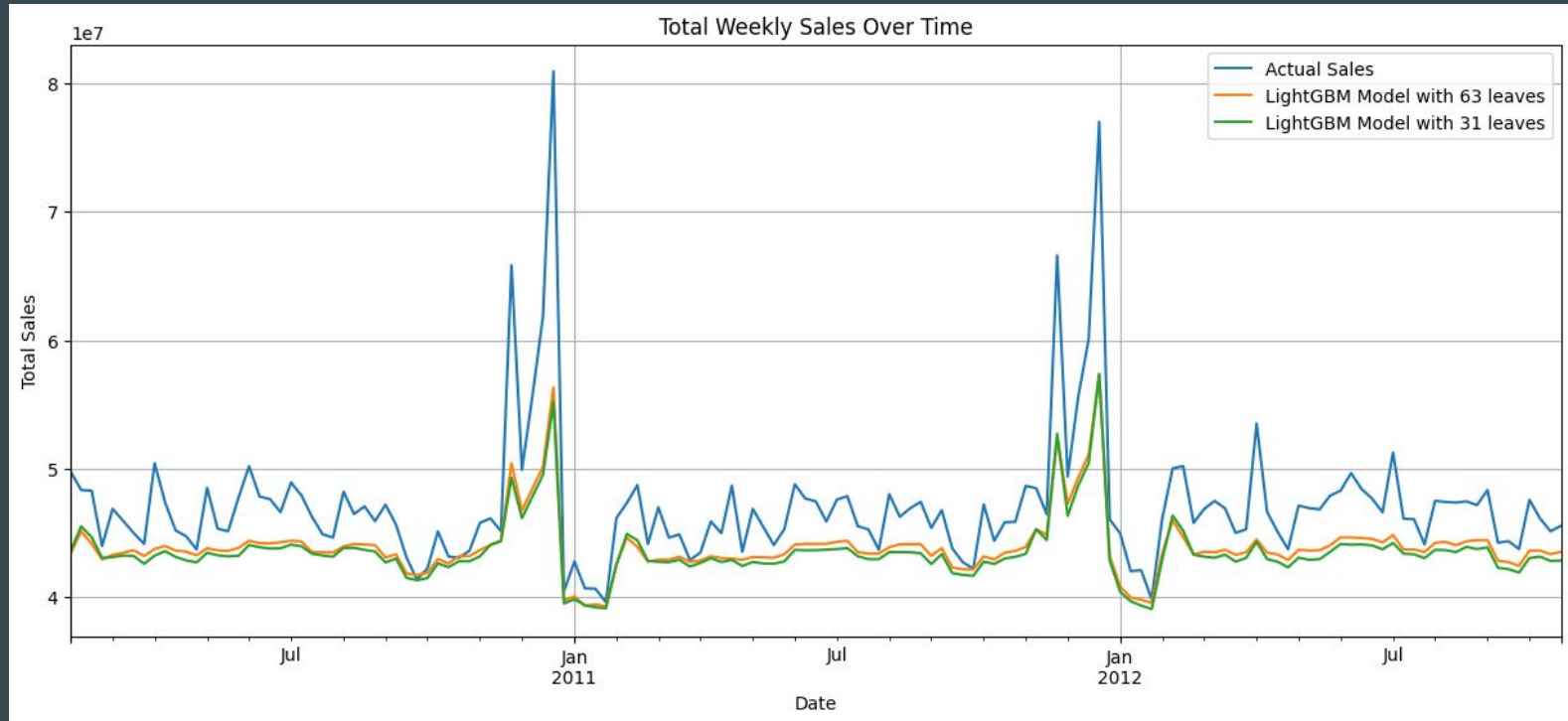
- Random Forests are not naturally suited for time series modeling because they do not account for the temporal order of data.
- While it's possible to add lag features or time-based variables through feature engineering, Random Forests treat observations as independent and identically distributed (i.i.d.), ignoring any sequential dependency between them.
- Moreover, since Random Forests randomly sample both features and data in parallel without preserving order, modeling temporal dynamics (like trends or seasonality) is very difficult. that is shown in the plots, even though its WMAE score is better than that of the xgboost, it can't predict trends such as holiday sales

LIGHTGBM

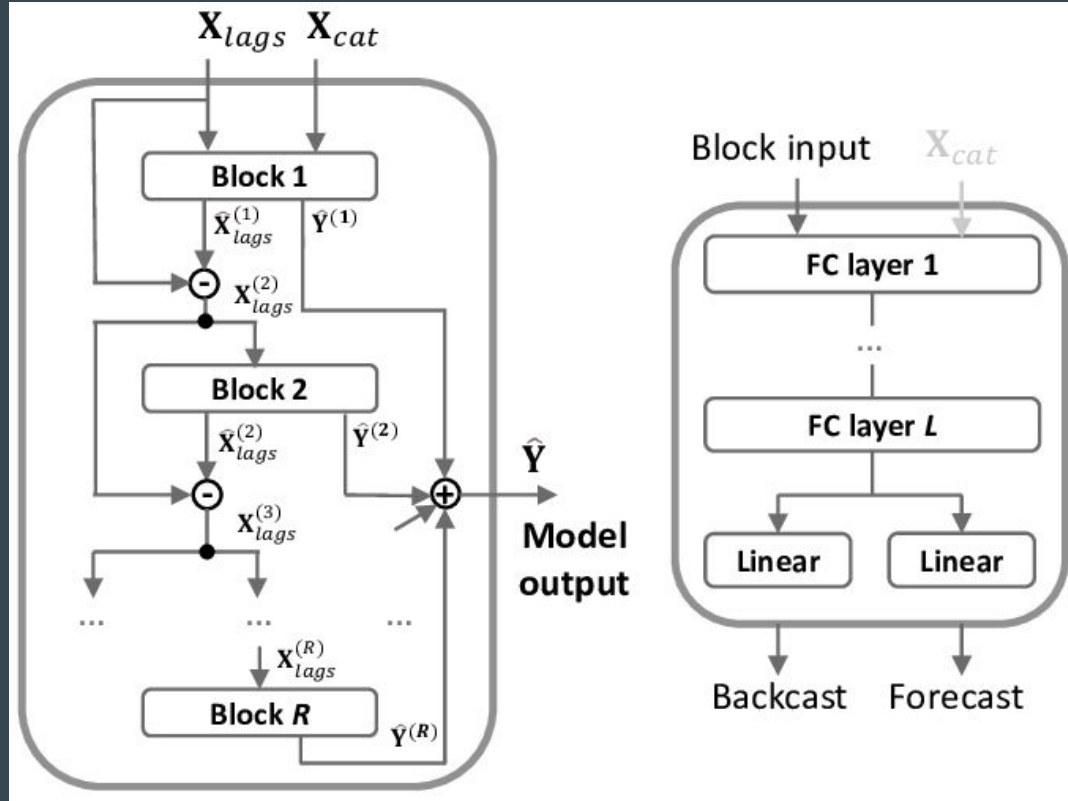
LightGBM was the best performing Decision Tree model for this problem, it had the highest train, validation and kaggle scores.



The Parameter which made the biggest improvement was the number of leaves, high amount of leaves (in our case, 63). number of estimators, reg_alpha and reg_lambda certainly affected the model, but not at the same magnitude

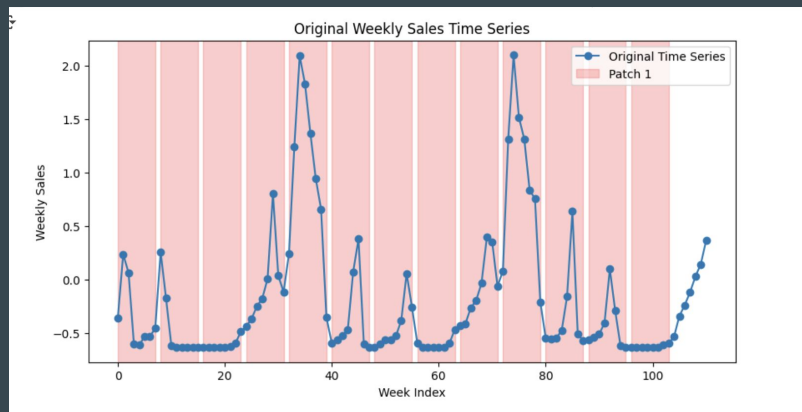


NBEATS



PatchTST

- **PatchTST (Patch-based Time Series Transformer)** is a deep learning model tailored for **univariate and multivariate time series forecasting**.
- Unlike traditional transformers, PatchTST divides the input series into **non-overlapping patches** and learns dependencies between them, making it both efficient and accurate.
- It avoids complex encoder-decoder setups by using a **simplified encoder-only architecture**.
- Handles **seasonality and long-range dependencies**
- Efficiently models **multivariate time series**, which suits our rich feature set: temperature, holidays, markdowns, CPI, unemployment, etc.
- Is highly scalable and works well with large datasets



Grid Search

```
sweep_config = {  
    'method': 'grid',  
    'name': 'patchtst_grid_search',  
    'metric': {'goal': 'minimize', 'name': 'val_wmae'},  
    'parameters': {  
        'seq_len': {'values': [26]},  
        'pred_len': {'values': [1]},  
        'patch_len': {'values': [4, 8, 12]},  
        'stride': {'values': [4, 8, 12]},  
        'd_model': {'values': [32, 64]},  
        'learning_rate': {'values': [1e-4, 1e-3]},  
        'n_layers': {'values': [1, 2]},  
        'epochs': {'values': [10]}  
    }  
}
```

seq_len: Number of past time steps fed into the model (input window).

pred_len: Number of future steps to predict

patch_len: Length of each input patch – smaller patches capture finer patterns, longer ones capture broader trends. |

Stride: How much to shift between patches – controls overlap and diversity.

D_model: Dimensionality of the model embeddings (internal feature space).

Learning_rate: Controls how quickly the model adapts weights during training.

N_layers: Depth of the transformer (more layers → more learning capacity, but risk of overfitting).

Epochs: Total training iterations through the dataset.

Best Configuration and Result

```
d_model: 32  
epochs: 10  
learning_rate: 0.0001  
patch_len: 12  
seq_len: 26  
stride: 4
```

Val WMAE: 2957.1

Final Results on Kaggle



submission (6).csv

Complete (after deadline) · 13h ago · XGBoost without Lag features

7116.07521

7014.83048



submission (5).csv

Complete (after deadline) · 19h ago · XGBoost with lag features and best parameters

6037.21085

5860.71966



submission (4).csv

Complete (after deadline) · 21h ago · Random Forest Model

7261.05578

7128.24506



submission (3).csv

Complete (after deadline) · 21h ago · Light GBM with low leaves

5773.85751

5650.46062



submission (2).csv

Complete (after deadline) · 1d ago · LightBGM model

5104.75333

4952.52175