# A survey of ML Based sales prediction

**Team**: Hinton's Heuristics

**Members**:
JatinKumar Janshali (24M2021)
Prashik Patil (24M2014)
Shoaib Ahamed (24M2102) [1]

**Course Instructor:**
Prof. Preethi Jyothi

# Task Description

| Aim: | • Predict Walmart's sales for the next **28 days** using the **M5** dataset. |
|---|---|
| Dataset: | • **6 years** of Walmart sales data across **3 states.** |
| Method: | • Incorporate events, holidays, and promotions into the forecast. |
| Focus: | • Evaluate **Linear Regression**, **LSTM** and **LightGBM** models. |
| Evaluation: | • The RMSSE formula from the M5 competition is the following: |

$$RMSSE = \sqrt{\frac{1}{h}\frac{\sum_{t=n+1}^{n+h}\left(Y_t - \widehat{Y}_t\right)^2}{\frac{1}{n-1}\sum_{t=2}^{n}(Y_t - Y_{t-1})^2}},$$

# Dataset Details

# Data Preprocessing

## Memory Optimization with Down casting

- **Purpose:** large data types (**float64**) into smaller ones (**float32**) w/o losing precision.

- **Implementation:** Dynamically using NumPy.

- **Impact:** Optimized data storage and faster processing

- **Result:** Reduction in size upto 45% *.

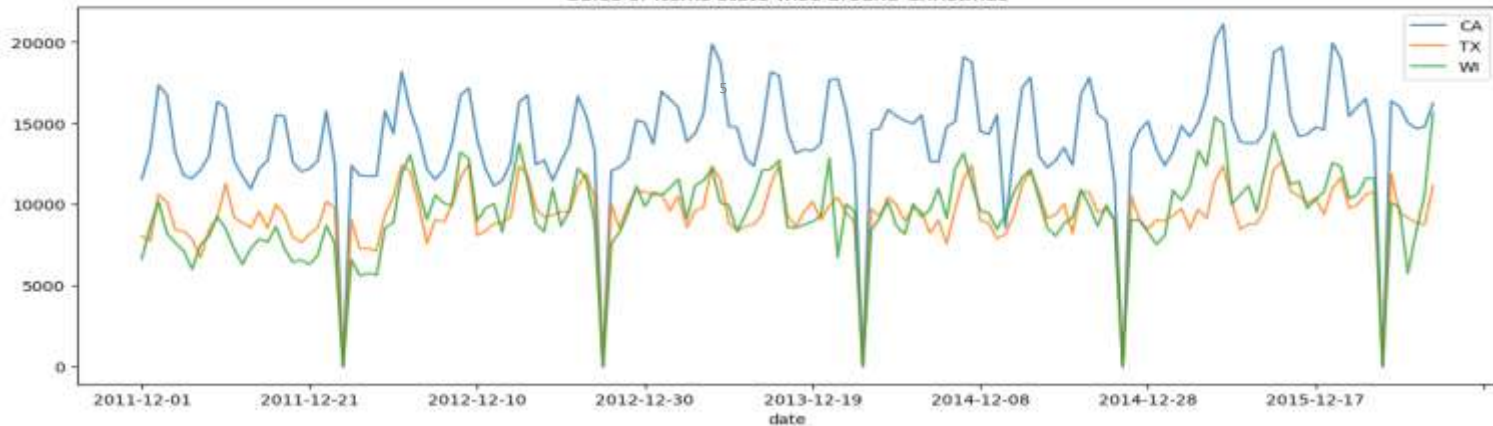## Data Normalization with Min-Max Scaling

- **Purpose:** Standardize numeric data to a common scale (-1 to 1 in this case).

- **Implementation:** Utilizes **MinMaxScaler** while preserving the proportionality in features

- **Impact:** training stability and convergence for features with varying scales.
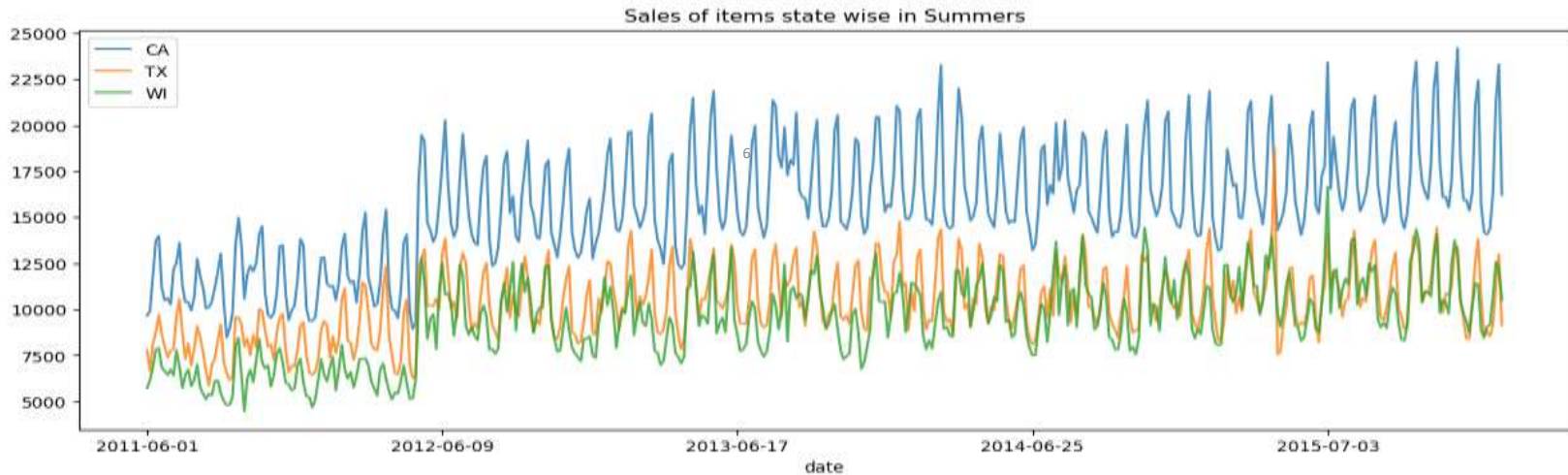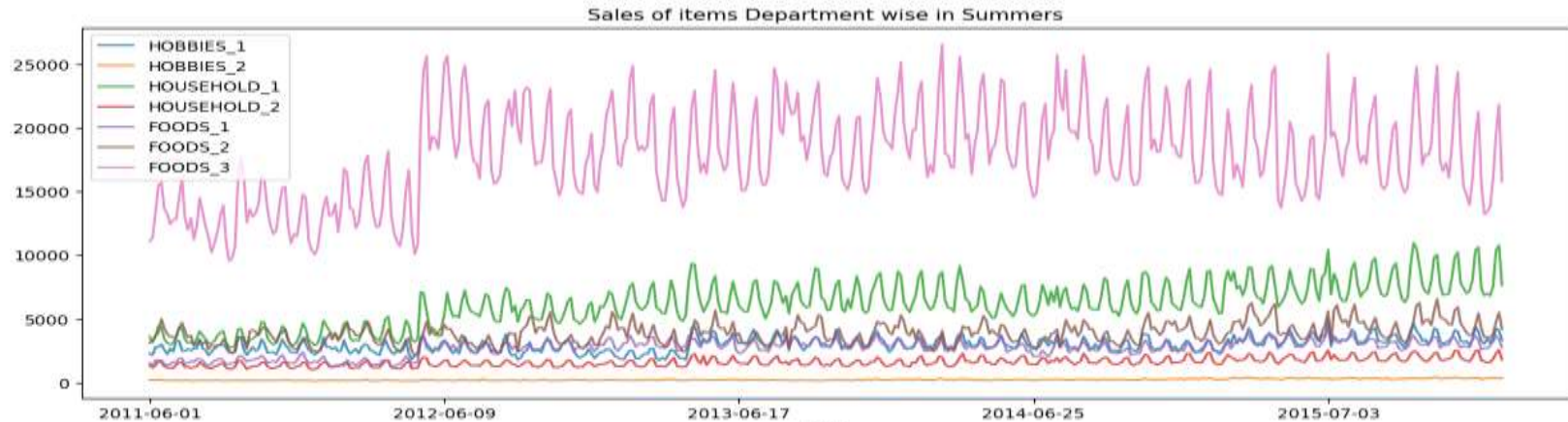
# EDA



Sales of items category wise during Christmas

Sales of items state wise around Christmas

# EDA



Sales of items Department wise in Summers

Sales of items state wise in Summers

# Linear Regression: Ridge regression model

**Salient Features:**

**Tailored Predictions**: Models can be optimized for specific category patterns and trends.

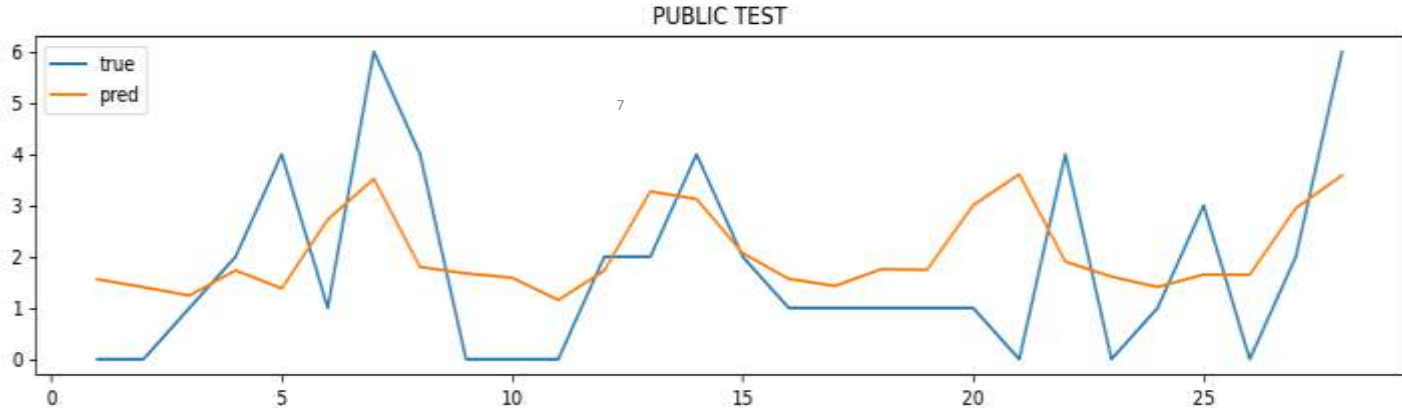**Reduced Noise**: Less diverse data leads to better focus on relevant features.

**Faster Training and Inference**: Smaller datasets accelerate model training and prediction.

**Improved Interpretability**: Insights from individual models can inform targeted strategies.

**Enhanced Accuracy**: Specialized models can yield more accurate predictions.

**RMSSE Score - 0.69725**          Sample Prediction for HOBBIES_1_004_CA_1



PUBLIC TEST

# LSTM

**Key Training Parameters**

Number of Epochs: 500      Number of Layers: 4
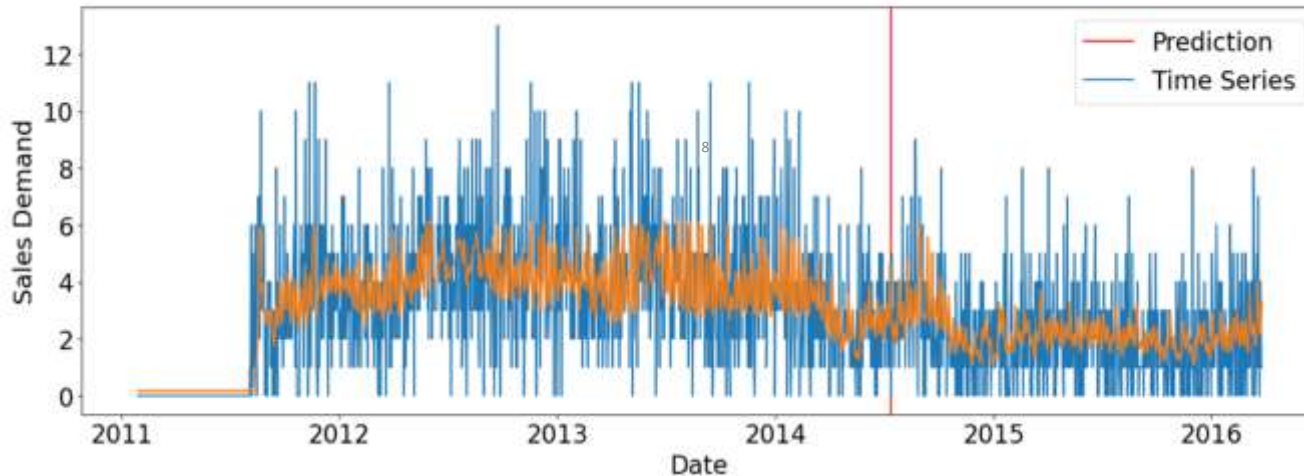
Learning Rate: $1\times10^{-3}$     Optimizer: Adam with L2 Regularization

Hidden Size: 512       Learning Rate Scheduler: patience = 50, factor = 0.5

**RMSSE - 0.76605**



Time-Series Prediction Entire Set

# LIGHTGBM

**Purpose**: To understand LightGBM and explore its features

**Approach Taken:**
Reduced the dataset from 6 years to 2 years
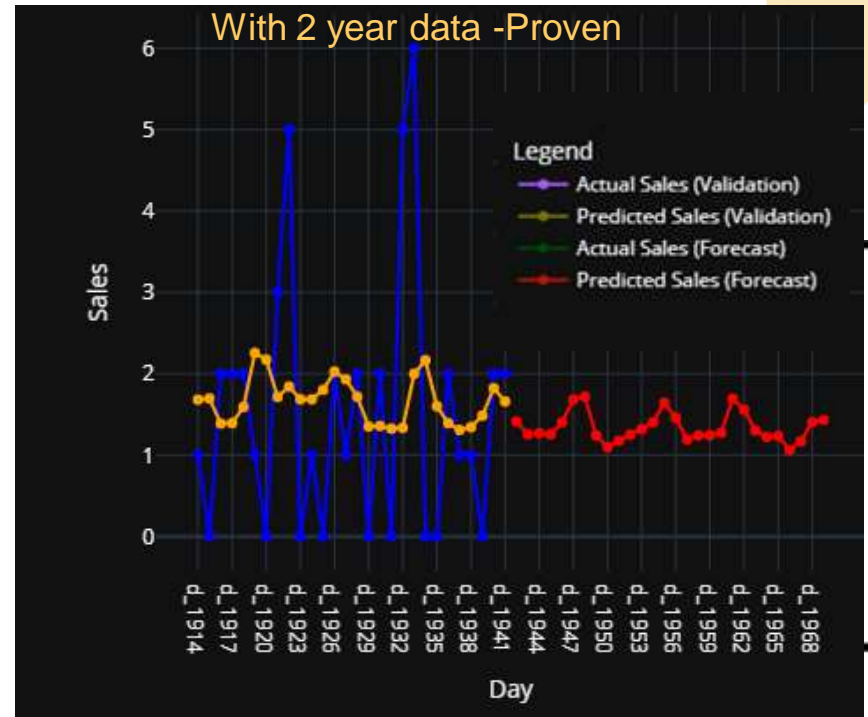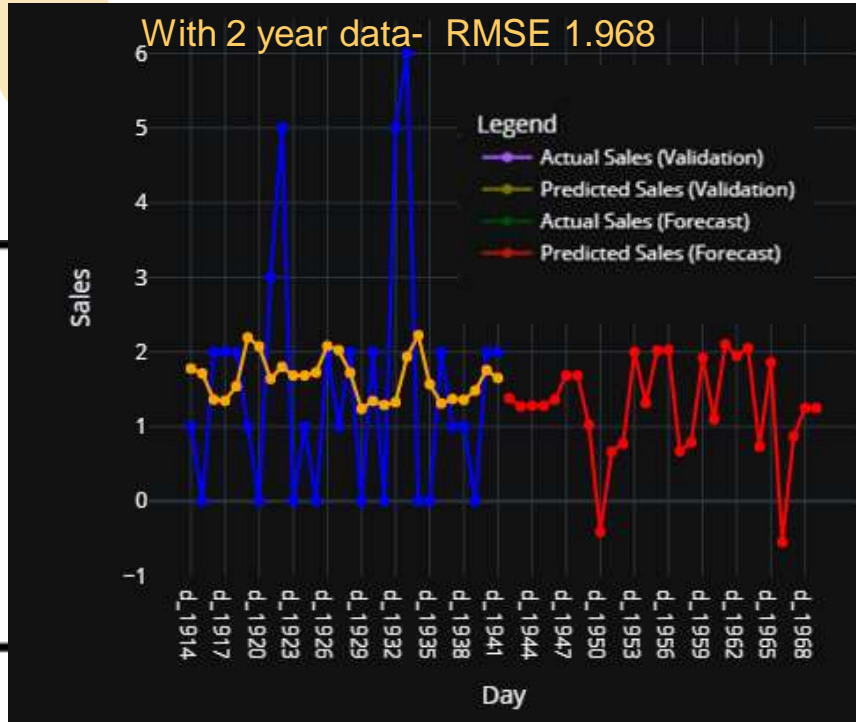Prediction carried out for random item selected from the dataset for next 28 days

**Feature Engineering:**
Merging relevant features of three dataset into one
Lag and rolling mean for last 7 days and 28 days

**Evaluation Methodology:**
Generated forecast for random items and compare it with proven prediction from kaggle

**Model Parameters**:

```
lgb_params = {
    'objective': 'regression',
    'metric': 'rmse',
    'boosting_type': 'gbdt',
    'num_leaves': 31,
    'learning_rate': 0.1,
    'feature_fraction': 0.5,
    'seed' : 2000
}
```

# LIGHTGBM-Results



**Possible Reasons for difference in results:**
1. Tweedie Regression objective function and other parameters considered for tuning
2. Recursive prediction

# TEAM TASK ALLOTMENT

| Sr. | Task | Executed By |
|---|---|---|
| 1 | Data Preprocessing | Prashik |
| 2 | EDA | Team |
| 3 | LSTM | Shoaib |
| 4 | Linear Regression | Prashik |
| 5 | LGBM | JatinKumar |

**Mr Ambuj Nayan(24M0003) has not participated in any project activity for unknown reasons.**