# Predict Future Sales
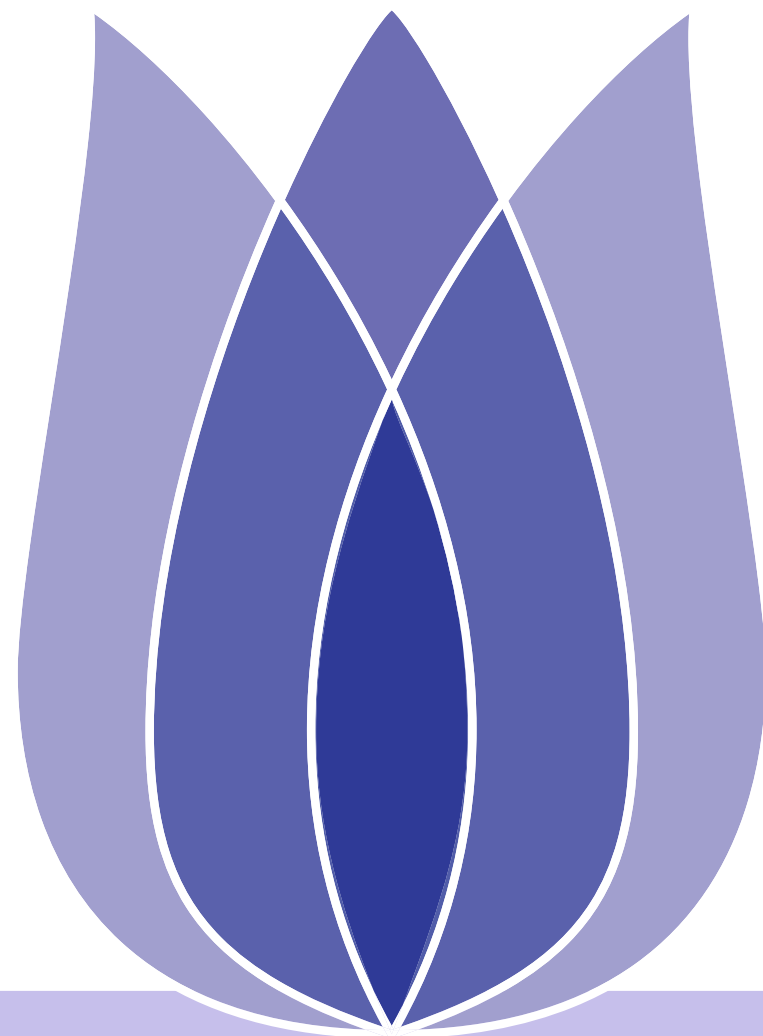
MingZhu Kang

Xi'an Shiyou University

Chinese Academy of Sciences

November 20, 2020

# Introduction

Defn

■ Project Introduction

Analyze the company's operating status, find out the relevant factors affecting the sales volume of goods, Taking the historical sales data set of convenience stores as the object of study, the data were preprocessed and feature extracted, and the model was used to train the data set to predict the sales volume of different goods in each store of the company in the next month.

TULIP *Team for Universal Learning and Intelligent Processing*

# Data Set Preprocessing

*TULIP Team for Universal Learning and Intelligent Processing*

# Preprocessing Of Project Data Sets

- Data Collection

  Download dataset from the kaggle project.

# Data Cleaning

# Training Set Data Cleaning

■ Use a scatter plot to observe the distribution of commodity prices and daily sales.

■ Filter for anomalies and apparent outliers



Figure 1: Distribution
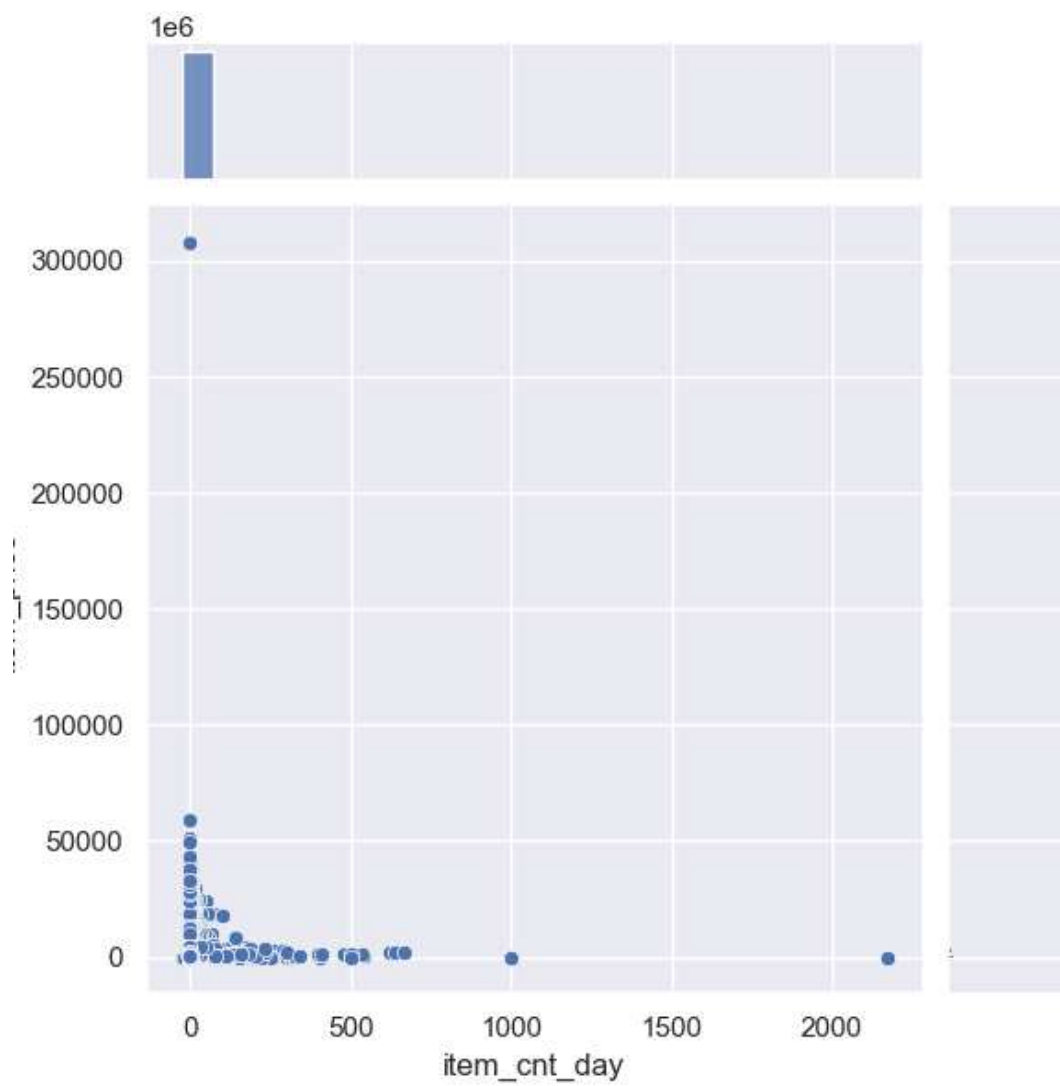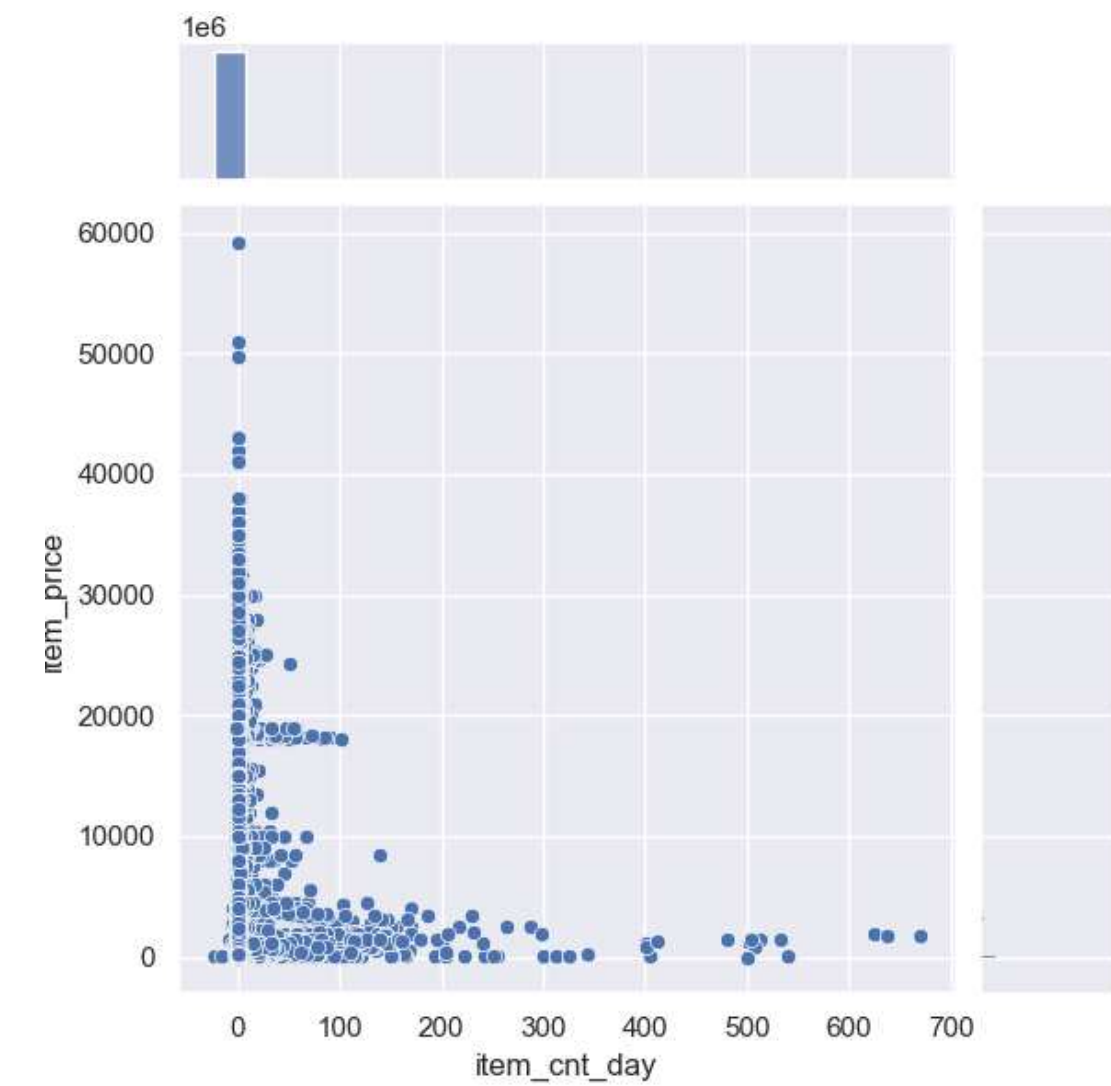


Figure 2: Filter Abnormal

TULIP *Team for Universal Learning and Intelligent Processing*

# Structured Data And Analysis

# Sales Analysis

Then, we created additional features. More specifically:

- Sales analysis

  Overall sales were down, and monthly sales were mostly down year on year. One item sold exceptionally well.
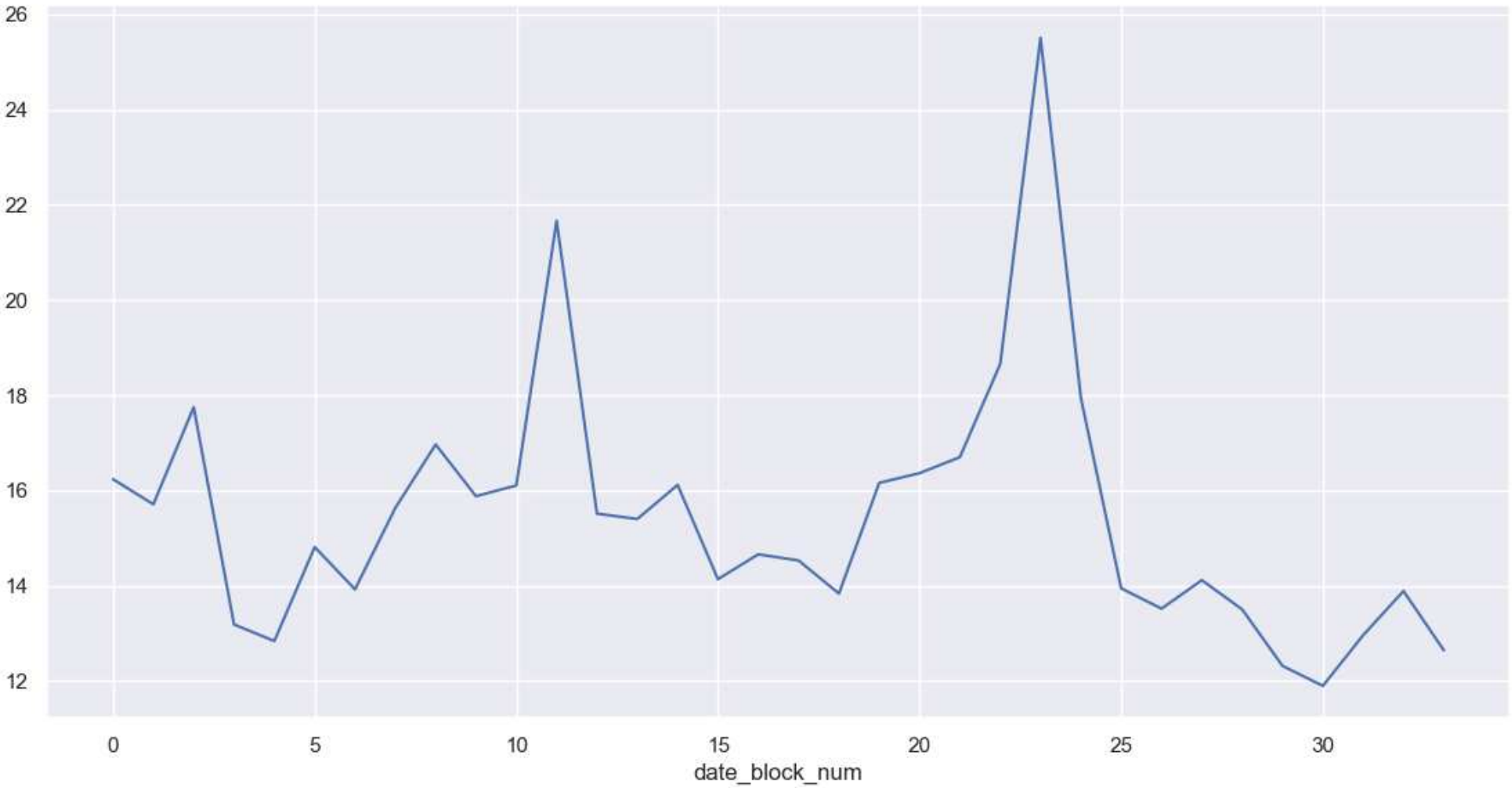
# Sales Analysis

- ■ Average number of items sold in the month
  In 2013 and 2014, the average monthly sales volume of goods under sale was basically 13-16, while in 2015, the average monthly sales volume of goods under sale decreased to 12-14.

TULIP *Team for Universal Learning and Intelligent Processing*

# Profit Analysis

# Profit Analysis

■ turnover analysis

# Profit Analysis

- The highest-grossing commodity

  The number one item in total revenue accounts for a percentage of monthly revenue

# Predict Future Trend

# Working With Training Set

- Handle closed stores and discontinued goods.
- Only keep the goods that are normally operated in the last 6 months and the goods with sale volume.



Figure 3: Closed Stores



Figure 4: Normal Opreation

■ use historical sales data to predict future sales.

Using the historical sales data as the characteristics of the model,

this month's sales results as labels to build a model for regression analysis.



Figure 5: Fusion Feature

# Model Adopted

# LightGBM Model

This project uses lightGBM model for training.

LightGBM is a fast, distributed, high-performance gradient enhancement framework based on decision tree algorithms.It supports category characteristics.

LightGBM supports category characteristics directly and natively by changing the decision rules of the decision tree algorithm, without transformation.
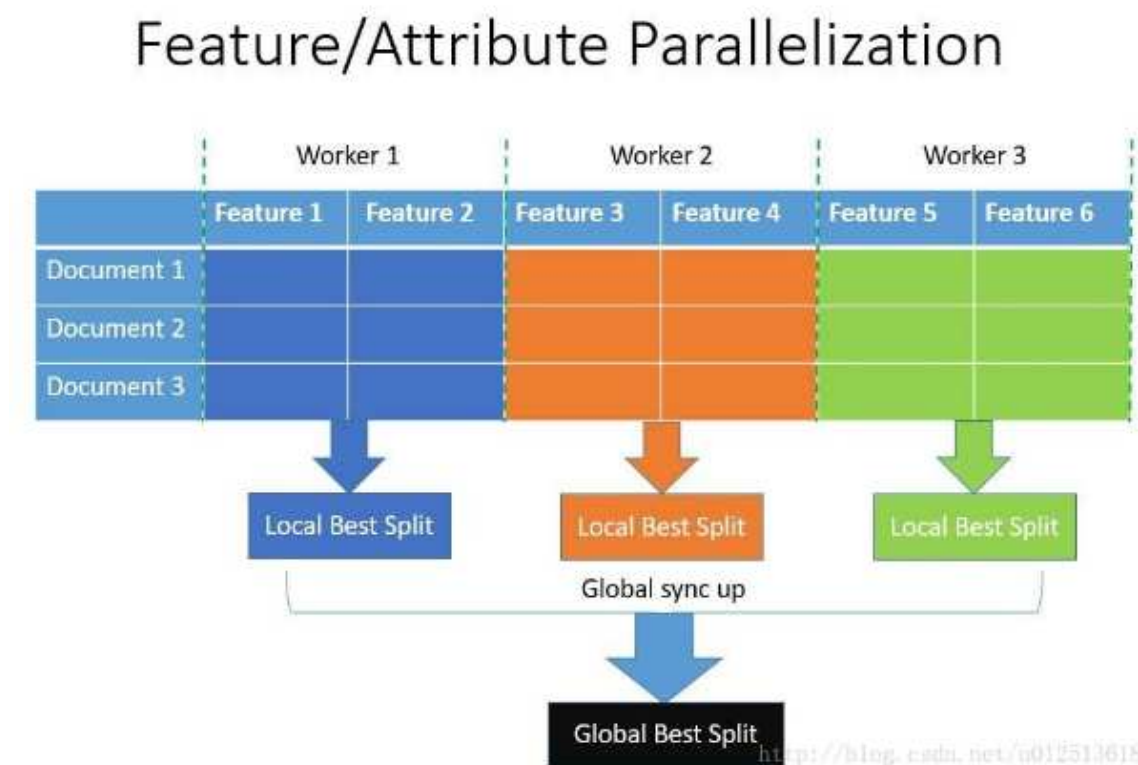


Figure 6: Feature Parallelization

Team for Universal Learning and Intelligent Processing

# Summary

# Project Summary

From data analysis methods to feature engineering and prediction model construction, a lot of time has been spent to study and comb.

Through this project, I have learned a lot, including the effective aspects of problem cutting, code implementation of analysis algorithm, design of analysis process, etc. which enables me to better grasp the thinking of data analysis on the whole.

In the process of predictive analysis, the theoretical and data support for feature analysis and model construction is not concise and powerful enough,which needs to be strengthened.