# Convenience Insights

**Multicampus Data Analysis & Engineer 34th Course**
**2024-03-07 ~ 2024-03-29**

# CONTENT

**INTRODUCTION**

- TEAM MEMBERS
- BACKGROUND
- WEB SERVICE
- DEVELOPMENT
  ENVIRONMENT
- PROJECT PERIOD

**PROJECT**

- FLOW CHART
- WBS

**DATA**

- DATA COLLECTION
- DATA PREPROCESSING

**EDA**

- EDA
- CORRELATION
  ANALYSIS
- MULTICOLLINEARITY

**MODELING**

- MODEL
- MODEL TRAINING
- MODEL PERFORMANCE
  VALIDATION
- CONVENIENCE STORE
  SALES PREDICTION

**APP & DOCUMENT**

- STREAMLIT
- LIMITATIONS
  IMPROVEMENTS
- REFERENCES
- APPENDIX

# CONTENT

**INTRODUCTION**

- TEAM MEMBERS
- BACKGROUND
- WEB SERVICE
- DEVELOPMENT
  ENVIRONMENT
- PROJECT PERIOD

**PROJECT**

- FLOW CHART
- WBS

**DATA**

- DATA COLLECTION
- DATA PREPROCESSING

**EDA**

- EDA
- CORRELATION
  ANALYSIS
- MULTICOLLINEARITY

**MODELING**

- MODEL
- MODEL TRAINING
- MODEL PERFORMANCE
  VALIDATION
- CONVENIENCE STORE
  SALES PREDICTION

**APP & DOCUMENT**

- STREAMLIT
- LIMITATIONS
  IMPROVEMENTS
- REFERENCES
- APPENDIX

# 1. INTRODUCTION

**TEAM MEMBERS**

**HYERIN CHOI**

- TEAM LEADER
- DATA PREPROCESSING
- EDA
- MAP VISUALIZATION
- STREAMLIT
  IMPLEMENTATION

**JINA KIM**

- DATA PREPROCESSING
- EDA
- STATISTICAL ANALYSIS
- MACHINE LEARNING

**MIN SONG**

- EDA
- STATISTICAL ANALYSIS
- MACHINE LEARNING

**JUNHO SONG**

- EDA
- STATISTICAL ANALYSIS
- MACHINE LEARNING
- MAP VISUALIZATION
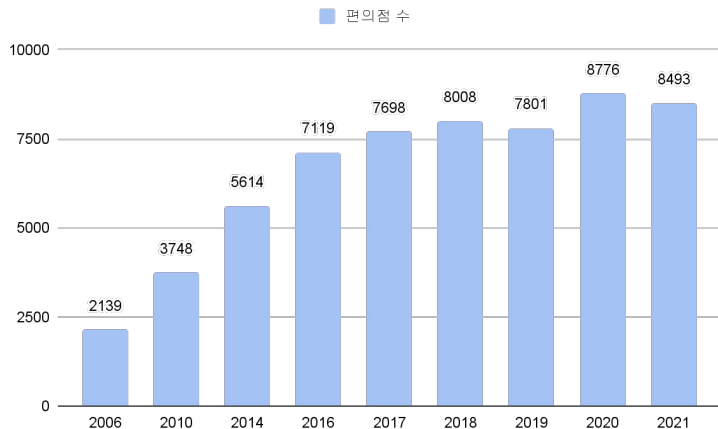- STREAMLIT
  IMPLEMENTATION

**DAEHEE HAN**

- STATISTICAL ANALYSIS
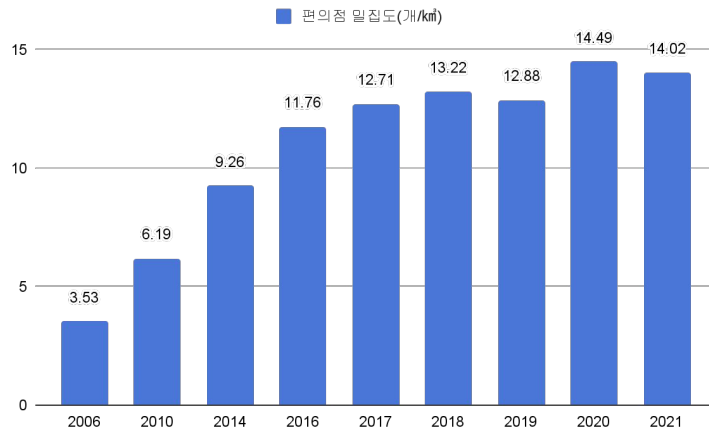- MACHINE LEARNING

# 1. INTRODUCTION

**BACKGROUND**

- As of the end of 2021, the number of convenience stores in Seoul totaled **8,493**, which is approximately **four times** the number compared to **2,139** in 2006.
- The convenience store density, measured by the number of stores per ㎢, also increased from 3.5 stores in 2006 to 14 stores in 2021.

| Convenience Store Count in Seoul |

■ 편의점 수



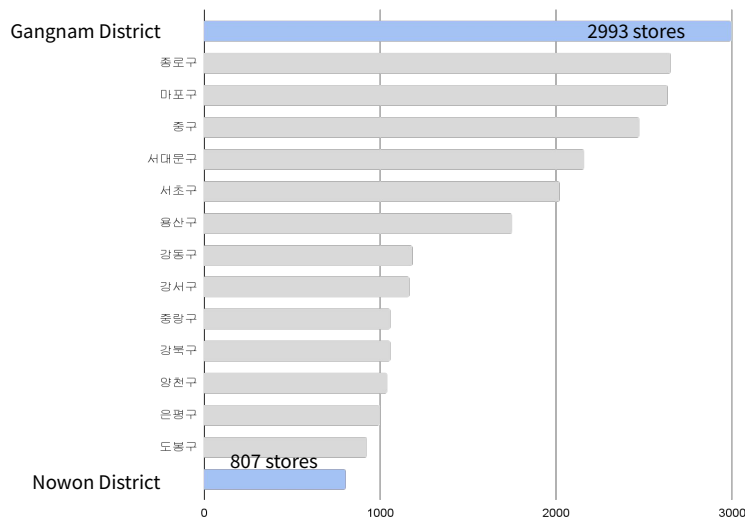| Convenience Store Density in Seoul |

■ 편의점 밀집도(개/㎢)



Source: Seoul Metropolitan Government, Operation Status and Current Situation Analysis Data of Convenience Stores in Seoul
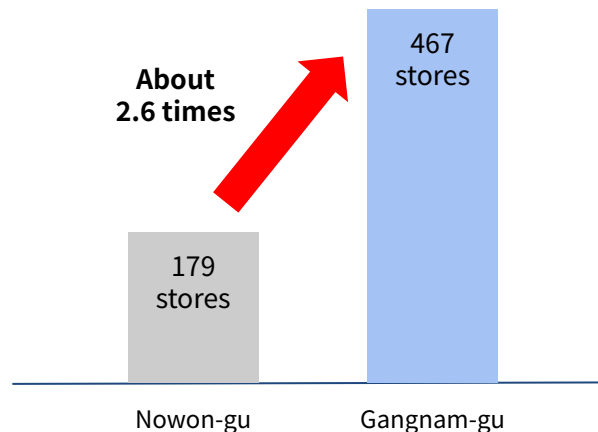
# 1. INTRODUCTION

**BACKGROUND**

- **Gangnam District** - 467 convenience stores, **Nowon District** - 179 convenience stores
  indicating a difference of approximately **2.6 times**.

- Suggesting that if one were to open a convenience store in Gangnam-gu, the competition would be more intense. Also implying
  the necessity of conducting a market analysis.

| Number of Convenience Facilities by District |

| Number of convenience stores in
Nowon-gu and Gangnam-gu |

# 1. INTRODUCTION

**BACKGROUND**

- When opening a convenience store franchise, there exists a **law that requires the franchisor to inform the franchisee about the expected sales figures**. However, **discrepancies between the projected and actual sales figures** can occur after the establishment, leading to difficulties for the franchisees.

**| The Problems in Opening Convenience Stores |**

⌂ 홈 > COVER STORY > 편의점의 역설

If the "projected sales figures" provided by the convenience store headquarters are exaggerated

HOME > 사회 > 사건/사고

Homeplus Inflating Projected Sales Figures for Prospective Convenience Store Owners

**[Convenience Store Dispute Cases]**

**Franchise Headquarters:** "Expected daily sales to reach approximately 1.3 million won"

**Franchisee A: "Actual daily sales only amount to 700,000 won"**

**[Planned Closure Turns Complicated...]**   [Source | Gyeonggi Fair Trade Support Center]

**Facing a penalty of 81 million won for contract termination**

**Facility and interior design costs borne by the headquarters**

**Franchisee, application for dispute resolution**

**Penalty reduced by 12 million won, Must pay 69 million won**

**[Hidden Issues]**

1 **Cases where expected sales are often verbally communicated**

2 **It's difficult to hold the headquarters accountable even if the expected sales differ**
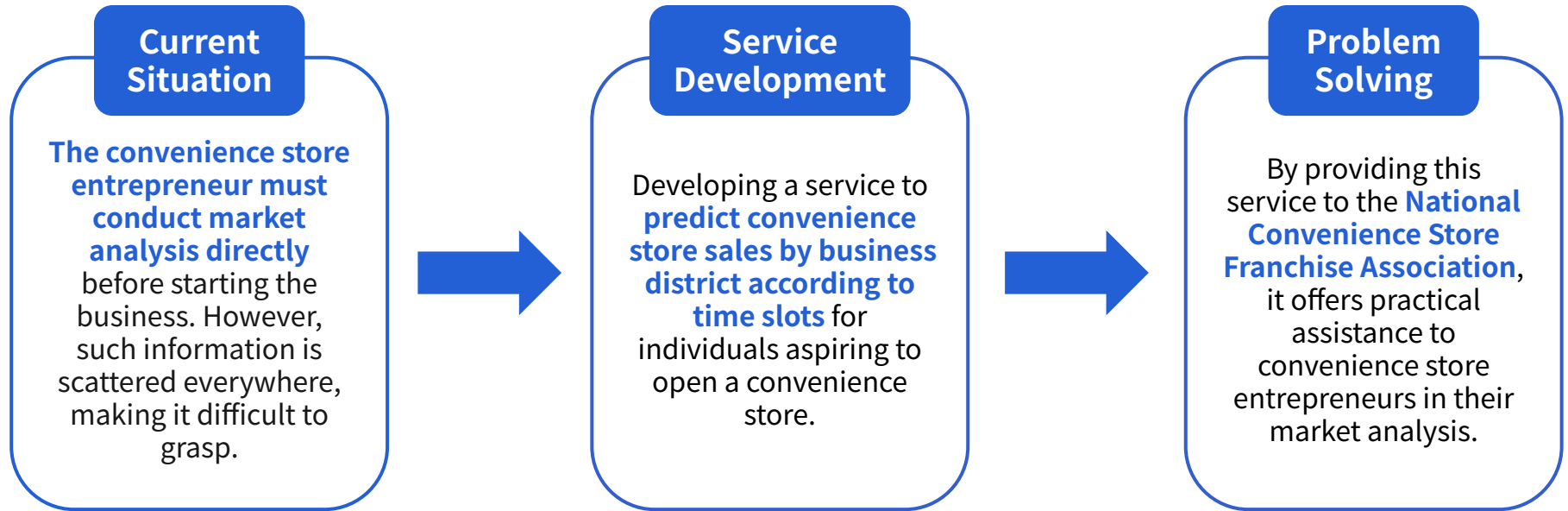
**Franchisees unable to close their business due to the burden of penalties**

Source: The Scoop, If the Projected Sales Figures Provided by the Convenience Store Headquarters Are Exaggerated…,
Citizen Daily, Homeplus Inflating Projected Sales Figures for Prospective Convenience Store Owners
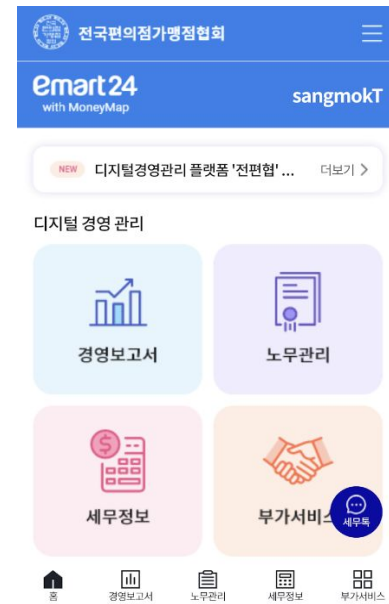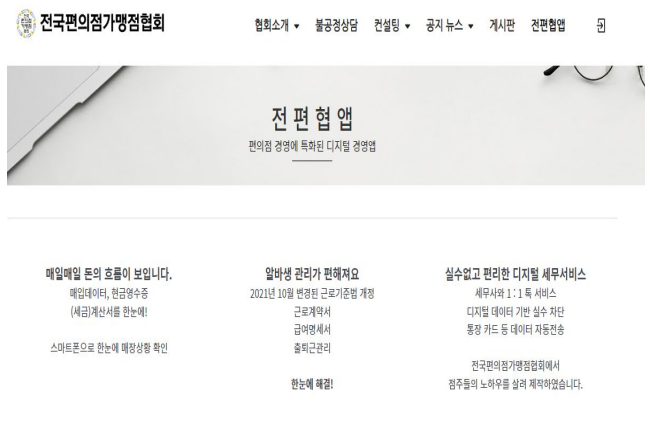
# 1. INTRODUCTION

**BACKGROUND**

## Current Situation

**The convenience store entrepreneur must conduct market analysis directly** before starting the business. However, such information is scattered everywhere, making it difficult to grasp.

## Service Development

Developing a service to **predict convenience store sales by business district according to time slots** for individuals aspiring to open a convenience store.

## Problem Solving

By providing this service to the **National Convenience Store Franchise Association**, it offers practical assistance to convenience store entrepreneurs in their market analysis.
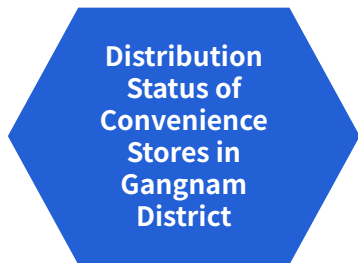
# 1. INTRODUCTION

**BACKGROUND**

- After accessing the website of the National Convenience Store Franchise Association, it has been confirmed that services such as management reporting, labor management, and tax information exist. However, there is no **market analysis service specifically designed for convenience store entrepreneurs**.

# 1. INTRODUCTION

**WEB SERVICE**

**Location by Business District and Average Sales**

**Distribution Status of Convenience Stores in Gangnam District**

**Convenience Store Sales Status in Gangnam District**

**Average Sales by Administrative District and Business District Time Slots**

**Top 5 Business Districts with Highest Sales by Time Slots**

**Sales Ranking**

**Sales Prediction**

**Expected Sales by Time Slot**

# 1. INTRODUCTION

**DEVELOPMENT ENVIRONMENT**

# 1. INTRODUCTION

**DEVELOPMENT ENVIRONMENT**

**Tool**

- Visual Studio
- Jupyter lab

**Deployment**

- Github
- Streamlit

**Language**

- Python

**Library**

- pandas
- numpy
- plotly
- matplotlib
- seaborn
- streamlt
- streamlit_option_menu
- streamlit_folium
- folium
- geopandas
- os
- time
- sklearn
- xgboost
- statsmodel
- scipy
- lgbm
- re
- jolib

# 1. INTRODUCTION

**PROJECT PERIOD**

| SUN | MON | TUE | WED | THU | FRI | SAT |
|---|---|---|---|---|---|---|
| 3 | 4 | 5 | 6 START | 7 | 8 | 9 |
| | | | SERVICE PLANNING | | | |
| | | | | DATA COLLECTION | | |
| | | | | | DATA PREPROCESSING | |
| 10 | 11 | 12 | 13 | 14 | 15 | 16 |
| DATA PREPROCESSING | | | | | | |
| EDA / STATISTICAL ANALYSIS | | | | | | |
| | | | | MODELING | | |
| 17 | 18 | 19 | 20 | 21 | 22 | 23 |
| MODELING | | | | | | |
| SERVICE IMPLEMENTATION | | | | | | |
| | | | | | PRESENTATION PREPARATION | |
| 24 | 25 | 26 | 27 | 28 FINAL | 29 FINISH | 30 |
| PRESENTATION PREPARATION | | | | | | |
| | | | | DEPLOYMENT | | |

# CONTENT

## INTRODUCTION



- TEAM MEMBERS
- BACKGROUND
- WEB SERVICE
- DEVELOPMENT
  ENVIRONMENT
- PROJECT PERIOD

## PROJECT



- FLOW CHART
- WBS

## DATA



- DATA COLLECTION
- DATA PREPROCESSING

## EDA



- EDA
- CORRELATION
  ANALYSIS
- MULTICOLLINEARITY

## MODELING



- MODEL
- MODEL TRAINING
- MODEL PERFORMANCE
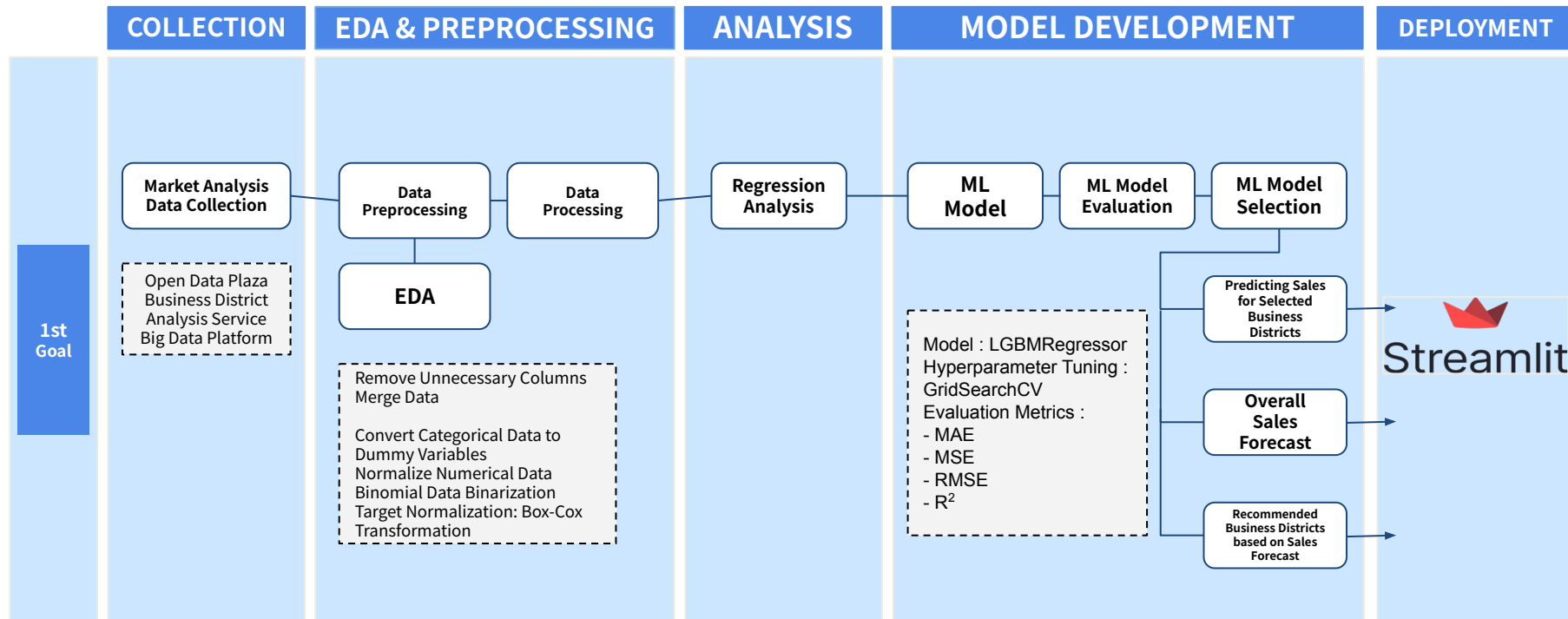  VALIDATION
- CONVENIENCE STORE
  SALES PREDICTION

## APP & DOCUMENT



- STREAMLIT
- LIMITATIONS
  IMPROVEMENTS
- REFERENCES
- APPENDIX

# 2. DOCUMENTATION

**FLOW CHART**

| COLLECTION | EDA & PREPROCESSING | ANALYSIS | MODEL DEVELOPMENT | DEPLOYMENT |
|---|---|---|---|---|

**1st Goal**

**Market Analysis Data Collection**

**Data Preprocessing** — **Data Processing**

**EDA**

Open Data Plaza
Business District
Analysis Service
Big Data Platform

Remove Unnecessary Columns
Merge Data

Convert Categorical Data to
Dummy Variables
Normalize Numerical Data
Binomial Data Binarization
Target Normalization: Box-Cox
Transformation

**Regression Analysis**

**ML Model** — **ML Model Evaluation** — **ML Model Selection**

Model : LGBMRegressor
Hyperparameter Tuning :
GridSearchCV
Evaluation Metrics :
- MAE
- MSE
- RMSE
- $R^2$

**Predicting Sales for Selected Business Districts**

**Overall Sales Forecast**

**Recommended Business Districts based on Sales Forecast**

Streamlit

# 2. DOCUMENTATION

## WBS

| CATEGORY | MAIN TASK | 1ST WEEK (3/06 ~ 3/09) | | | 2ND WEEK (3/10 ~ 3/16) | | | 3RD WEEK (3/17 ~ 3/23) | | | 4TH WEEK (3/24 ~ 3/29) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PLANNING | TOPIC SELECTION | ■ | ■ | ■ | | | | | | | | | |
| DATA COLLECTION | API DATA COLLECTION | | ■ | ■ | ■ | | | | | | | | |
| DATA PREPROCESSING | DATA FORMAT STANDARDIZATION | | | ■ | ■ | ■ | | | | | | | |
| DATA CHECK | EDA | | | | ■ | ■ | ■ | | | | | | |
| VISUALIZATION | SPACE / MAP VISUALIZATION | | | | | ■ | ■ | ■ | ■ | | | | |
| STATISTICAL ANALYSIS | REGRESSION ANALYSIS | | | | ■ | ■ | ■ | ■ | | | | | |
| MODELING | MODEL GENERATION / EVALUATION | | | | | ■ | ■ | ■ | ■ | | | | |
| SERVICE IMPLEMENTATION | STREAMLIT DEPLOYMENT | | | | | ■ | ■ | ■ | ■ | ■ | ■ | | |
| FINAL | FINAL PRESENTATION | | | | | | | | | | | | ■ |

# CONTENT

## INTRODUCTION

- TEAM MEMBERS
- BACKGROUND
- WEB SERVICE
- DEVELOPMENT
  ENVIRONMENT
- PROJECT PERIOD

## PROJECT

- FLOW CHART
- WBS

## DATA EXPLORATION

- DATA COLLECTION
- DATA PREPROCESSING

## EDA

- EDA
- CORRELATION
  ANALYSIS
- MULTICOLLINEARITY

## MODELING

- MODEL
- MODEL TRAINING
- MODEL PERFORMANCE
  VALIDATION
- CONVENIENCE STORE
  SALES PREDICTION

## APP & DOCUMENT

- STREAMLIT
- LIMITATIONS
  IMPROVEMENTS
- REFERENCES
- APPENDIX

# 3. DATA EXPLORATION

## DATA COLLECTION

| Major Category | Commercial Area Analysis Service Data | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Classification | Estimated Sales | Population by Road | Resident Population | Working Population | Income and Expenditure | Areas | Customer Attraction Facilities | Stores |
| Details | - Quarterly sales amount<br>- Number of sales transactions per quarter<br>- Sales amount by time slot<br>- Number of sales transactions by time slot<br>- Sales amount by age group<br>- Number of sales transactions by age group<br>- Sales amount by gender<br>- Number of sales transactions by gender<br>- Sales amount by day of the week<br>- Number of sales transactions by day of the week | - Total floating population<br>- Floating population by gender<br>- Floating population by age<br>- Floating population by time slot<br>- Floating population by day of the week | - Total resident population<br>- Resident population by gender<br>- Resident population by age group<br>- Total number of households<br>- Number of apartment households<br>- Number of non-apartment households | - Total working population<br>- Working population by gender<br>- Working population by age group | - Average monthly income<br>- Income level<br>- Total expenditure<br>- Type of expenditure | - X coordinate<br>- Y coordinate<br>- District code<br>- Neighborhood code<br>- Area size | - Total number of customer attraction facilities<br>- Number of customer attraction facilities by type | - Store opening/closing rate<br>- Number of stores opened/closed<br>- Total number of stores<br>- Number of stores by similar business type<br>- Number of franchise stores |
| Source | Seoul Open Data Plaza | | | | | | | |

# 3. DATA EXPLORATION

**DATA PREPROCESSING**

## 1. Extract only convenience store data in Gangnam District



서울시 상권분석서비스(추정매출-상권)

서울시 상권 영역 내의 점포들의 추정 매출 정보를 제공합니다.
※단위 : 원

일반행정

Seoul

Extract data where "Service_Industry_Code_Name" is convenience store



서울시 상권분석서비스(영역-상권)

서울시 상권분석 서비스에서 사용중인 상권영역 정보입니다.
※ 행정동 코드는 행정안전부에서 고시한 "주민등록 행정기관코드"를 사용하고 있습니다.
※좌표계 EPSG: 5181

일반행정

Extract data where "District_Code_Name" is convenience store.

# 3. DATA EXPLORATION

**DATA PREPROCESSING**

## 2. Checking for Missing Values Before Data Merge

| Category | Estimated Sales | Population by Road | Resident Population | Working Population | Income and Expenditure | Areas | Customer Attraction Facilities | Stores |
|---|---|---|---|---|---|---|---|---|
| **Null Values** | No Null Values | No Null Values | No Null Values | No Null Values | - Monthly Average Income Amount<br>- Income Bracket Code<br>- Total Expenditure Amount<br>- Total Grocery Expenditure Amount<br>- Total Clothing and Footwear Expenditure Amount<br>- Total Household Goods Expenditure Amount<br>- Total Medical Expenses Amount<br>- Total Transportation Expenditure Amount<br>- Total Leisure Expenditure Amount<br>- Total Cultural Expenditure Amount<br>- Total Education Expenditure Amount<br>- Total Entertainment Expenditure Amount | No Null Values | - Number of Government Offices<br>- Number of Banks<br>- Number of General Hospitals<br>- Number of Hospitals<br>- Number of Pharmacies<br>- Number of Kindergartens<br>- Number of Elementary Schools<br>- Number of High Schools<br>- Number of Universities<br>- Number of Department Stores<br>- Number of Supermarkets<br>- Number of Theaters<br>- Number of Accommodation Facilities<br>- Number of Airports<br>- Number of Railway Stations<br>- Number of Bus Terminals<br>- Number of Subway Stations<br>- Number of Bus Stops | No Null Values |

# 3. DATA EXPLORATION

**DATA PREPROCESSING**

## 2. Checking for Missing Values Before Data Merge

Replace with the average value of a different quarter in the same commercial district

| Category | Estimated Sales | Population by Road | Resident Population | Working Population | Income and Expenditure | Areas | Customer Attraction Facilities | Stores |
|---|---|---|---|---|---|---|---|---|
| Null Values | No Null Values | No Null Values | No Null Values | No Null Values | - Monthly Average Income Amount<br>- Income Bracket Code<br>- Total Expenditure Amount<br>- Total Grocery Expenditure Amount<br>- Total Clothing and Footwear Expenditure Amount<br>- Total Household Goods Expenditure Amount<br>- Total Medical Expenses Amount<br>- Total Transportation Expenditure Amount<br>- Total Leisure Expenditure Amount<br>- Total Cultural Expenditure Amount<br>- Total Education Expenditure Amount<br>- Total Entertainment Expenditure Amount | No Null Values | - Number of Government Offices<br>- Number of Banks<br>- Number of General Hospitals<br>- Number of Hospitals<br>- Number of Pharmacies<br>- Number of Kindergartens<br>- Number of Elementary Schools<br>- Number of High Schools<br>- Number of Universities<br>- Number of Department Stores<br>- Number of Supermarkets<br>- Number of Theaters<br>- Number of Accommodation Facilities<br>- Number of Airports<br>- Number of Railway Stations<br>- Number of Bus Terminals<br>- Number of Subway Stations<br>- Number of Bus Stops | No Null Values |

# 3. DATA EXPLORATION

**DATA PREPROCESSING**

## 2. Checking for Missing Values Before Data Merge

Replace Null Values with 0

| Category | Estimated Sales | Population by Road | Resident Population | Working Population | Income and Expenditure | Areas | Customer Attraction Facilities | Stores |
|---|---|---|---|---|---|---|---|---|
| Null Values | No Null Values | No Null Values | No Null Values | No Null Values | - Monthly Average Income Amount<br>- Income Bracket Code<br>- Total Expenditure Amount<br>- Total Grocery Expenditure Amount<br>- Total Clothing and Footwear Expenditure Amount<br>- Total Household Goods Expenditure Amount<br>- Total Medical Expenses Amount<br>- Total Transportation Expenditure Amount<br>- Total Leisure Expenditure Amount<br>- Total Cultural Expenditure Amount<br>- Total Education Expenditure Amount<br>- Total Entertainment Expenditure Amount | No Null Values | - Number of Government Offices<br>- Number of Banks<br>- Number of General Hospitals<br>- Number of Hospitals<br>- Number of Pharmacies<br>- Number of Kindergartens<br>- Number of Elementary Schools<br>- Number of High Schools<br>- Number of Universities<br>- Number of Department Stores<br>- Number of Supermarkets<br>- Number of Theaters<br>- Number of Accommodation Facilities<br>- Number of Airports<br>- Number of Railway Stations<br>- Number of Bus Terminals<br>- Number of Subway Stations<br>- Number of Bus Stops | No Null Values |

# 3. DATA EXPLORATION

## 3. Data Merge

| Category | Common Column |
|---|---|
| Estimated Sales | - Base Year and Quarter Code<br>- Commercial District Division Code<br>- Commercial District Division Code Name<br>- Commercial District Code<br>- Commercial District Code Name |
| Population by Road | |
| Resident Population | |
| Working Population | |
| Income and Expenditure | |
| Customer Attraction Facilities | |
| Stores | |

1st merge based on common columns

⬇

2nd merge with Gangnam District location information

# 3. DATA EXPLORATION

**DATA PREPROCESSING**

## 4. Changing the Columns

| Existing Column |
| :---: |
| Base Year and Quarter Code |

➡️

| Changed Column |
| :---: |
| Base_Year<br>Base_Quarter |

| Base_Year | Base_Quarter |
| :---: | :---: |
| 2021 | 1 |
| 2021 | 1 |
| 2021 | 1 |
| 2021 | 1 |
| 2021 | 1 |

| Existing Column |
| :---: |
| Sales Amount from xx to xx Hours |

➡️

| Changed Column |
| :---: |
| Time Period<br>Sales Amount by Time Period |

| Time Period | Sales Amount by Time Period |
| :---: | :---: |
| 00~06 | 377166450.0 |
| 06~11 | 222467605.0 |
| 11~14 | 192457360.0 |
| 14~17 | 230188421.0 |
| 17~21 | 531497598.0 |

# 3. DATA EXPLORATION

**DATA PREPROCESSING**

## 5. Removing Unnecessary Columns

| Unnecessary Columns |
| --- |
| - Sales Transaction<br>- Count Related Columns<br>- Gender Related Columns<br>- Day of the Week Related Columns<br>...<br>- Number of Visitor Facilities by Type<br>- Total Expenditure Amount by Type |

Remove Unnecessary Columns

| Final Columns | |
| --- | --- |
| - Sales Amount by Time Period<br>- Base Year<br>- Base Quarter<br>- Commercial District Division Name<br>- Commercial District Name<br>- Administrative Neighborhood Name<br>- Time Period<br>- Floating Population Count by Time Period<br>- Total Working Population<br>- Total Resident Population<br>- Total Number of Households<br>- Number of Attraction Facilities | - Monthly Average Income Amount<br>- Total Expenditure Amount<br>- Number of Stores in Similar Industries<br>- Number of Newly Opened Stores<br>- Number of Closed Stores |

# CONTENT

## INTRODUCTION

- TEAM MEMBERS
- BACKGROUND
- WEB SERVICE
- DEVELOPMENT
  ENVIRONMENT
- PROJECT PERIOD

## PROJECT

- FLOW CHART
- WBS

## DATA

- DATA COLLECTION
- DATA PREPROCESSING

## EDA

- EDA
- CORRELATION
  ANALYSIS
- MULTICOLLINEARITY

## MODELING

- MODEL
- MODEL TRAINING
- MODEL PERFORMANCE
  VALIDATION
- CONVENIENCE STORE
  SALES PREDICTION

## APP & DOCUMENT

- STREAMLIT
- LIMITATIONS
  IMPROVEMENTS
- REFERENCES
- APPENDIX

# 4. EDA

**Dependent Variable**

## Final Columns

- Sales Amount by Time Period
- Base Year
- Base Quarter
- Commercial District Division Name
- Commercial District Name
- Administrative Neighborhood Name
- Time Period
- Floating Population Count by Time Period
- Total Working Population
- Total Resident Population
- Total Number of Households
- Number of Attraction Facilities

- Monthly Average Income Amount
- Total Expenditure Amount
- Number of Stores in Similar Industries
- Number of Newly Opened Stores
- Number of Closed Stores

The EDA for each variable can be found in the appendix.

# 4. EDA

**FINAL COLUMNS CHECK**

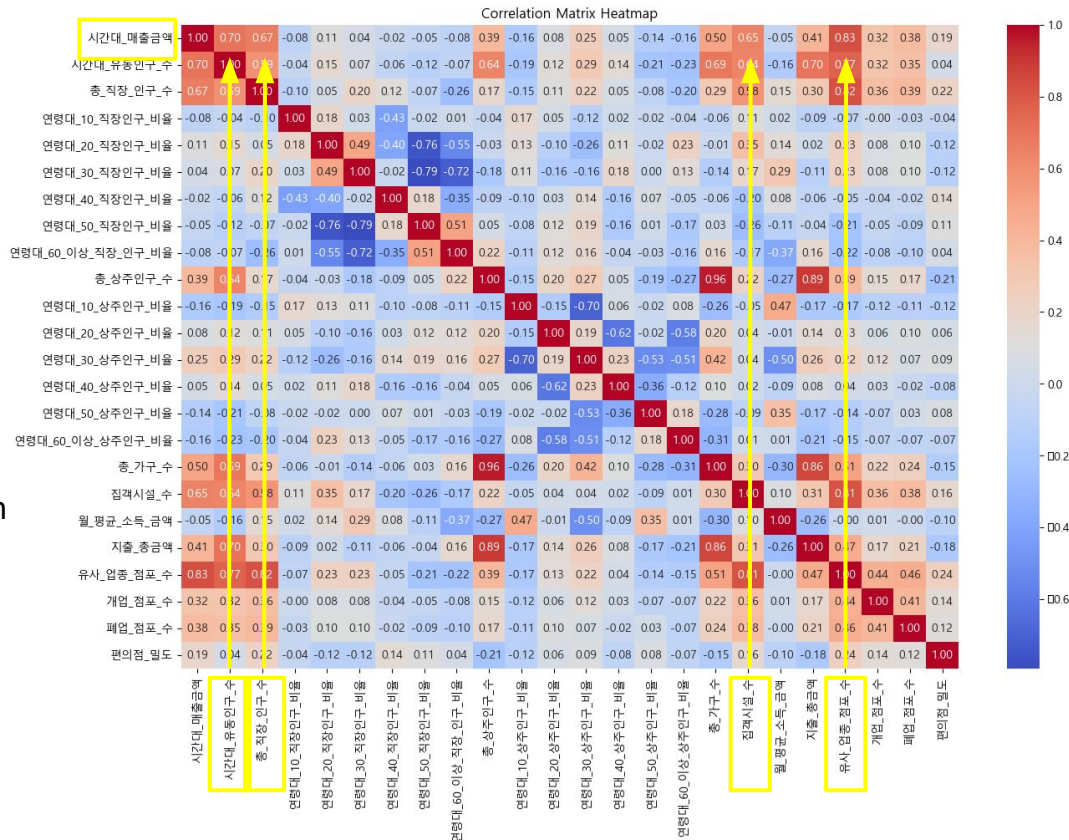| Final Columns | | |
|---|---|---|
| - Sales Amount by Time Period<br>- Base Year<br>- Base Quarter<br>- Commercial District Division Name<br>- Commercial District Name<br>- Administrative Neighborhood Name<br>- Time Period<br>- Floating Population Count by Time Period<br>- Total Working Population<br>- Total Resident Population<br>- Total Number of Households<br>- Number of Attraction Facilities | - Monthly Average Income Amount<br>- Total Expenditure Amount<br>- Number of Stores in Similar Industries<br>- Number of Newly Opened Stores<br>- Number of Closed Stores | **Independent Variables** |

The EDA for each variable can be found in the appendix.

# 4. EDA

**EDA**

- In Gangnam District, the number of convenience stores closing is generally higher than the number of new openings.
- Although the number of convenience stores is slightly decreasing, it is recording levels that are almost the same.
- On the other hand, sales amounts are showing a trend of gradual increase.

| Quarterly Convenience Store Numbers and Openings/Closures |



| Year and Quarterly Sales Amount |

# 4. EDA

**EDA**

- Sales vary by commercial district.
- When extracting a random commercial district, the sales trends by time of day are similar, but each district shows slightly different sales trends.
- It can be expected that the factors affecting sales by time of day vary depending on the commercial district.

| Average Sales by Commercial District from 2021 to Q3 2023 |



| Average Sales per Store by Time of Day According to Commercial District |

# 4. EDA

**EDA**

- High sales are recorded during commuting hours.
- Sales tend to increase when there is a high foot traffic, but despite a significant decrease in foot traffic from 21:00 to 24:00, sales during this time are the next highest after commuting hours.
- This suggests that, in addition to the number of pedestrians, other variables are also expected to influence sales by time of day.



| Average Sales Amount and Average Pedestrian Traffic by Time of Day |

# 4. EDA

**ADDING COLUMNS**

- Based on prior research, add the columns 'Convenience Store Density', 'Proportion of Resident Population by Age Group', and 'Proportion of Working Population by Age Group'.

## Existing Columns

| | |
|---|---|
| - Sales Amount by Time Period | - Monthly Average Income Amount |
| - Base Year | |
| - Base Quarter | - Total Expenditure Amount |
| - Commercial District Division Name | - Number of Stores in Similar Industries |
| - Commercial District Name | - Number of Newly Opened Stores |
| - Administrative Neighborhood Name | |
| - Time Period | - Number of Closed Stores |
| - Floating Population Count by Time Period | |
| - Total Working Population | |
| - Total Resident Population | |
| - Total Number of Households | |
| - Number of Attraction Facilities | |

## Final Columns

| | |
|---|---|
| - Sales Amount by Time Period | - Monthly Average Income Amount |
| - Base Year | |
| - Base Quarter | - Total Expenditure Amount |
| - Commercial District Division Name | - Number of Stores in Similar Industries |
| - Commercial District Name | - Number of Newly Opened Stores |
| - Administrative Neighborhood Name | |
| - Time Period | - Number of Closed Stores |
| - Floating Population Count by Time Period | - Convenience Store Density |
| - Total Working Population | - Proportion of Resident Population by Age |
| - Total Resident Population | - Proportion of Working Population by Age |
| - Total Number of Households | |
| - Number of Attraction Facilities | |

# 4. EDA

**CORRELATION ANALYSIS**

- Number of Stores in Similar Industries - 0.83
- Floating Population Count by Time Period - 0.70
- Total Working Population - 0.67
- Number of Attraction Facilities - 0.65

=> are expected to have a relatively significant impact on "Sales Amount by Time Segment".



Correlation Matrix Heatmap

# 4. EDA

**MULTICOLLINEARITY**

- Checking the Variance Inflation Factor (VIF) for the independent variables
- Proportion of Working Population by Age, Proportion of Resident Population by Age, Total Resident Population, Total Households, and Number of Similar Business Stores—is high, exceeding 10.

=> This suggests the presence of multicollinearity among these variables. It is expected that utilizing tree-based models will help to address the issue of multicollinearity.

| | VIF Factor | features |
|---|---|---|
| 0 | 4.457923 | 시간대_유동인구_수 |
| 1 | 5.303403 | 총_직장_인구_수 |
| 2 | inf | 연령대_10_직장인구_비율 |
| 3 | inf | 연령대_20_직장인구_비율 |
| 4 | inf | 연령대_30_직장인구_비율 |
| 5 | inf | 연령대_40_직장인구_비율 |
| 6 | inf | 연령대_50_직장인구_비율 |
| 7 | inf | 연령대_60_이상_직장_인구_비율 |
| 8 | 35.345023 | 총_상주인구_수 |

| | | |
|---|---|---|
| 9 | inf | 연령대_10_상주인구_비율 |
| 10 | inf | 연령대_20_상주인구_비율 |
| 11 | inf | 연령대_30_상주인구_비율 |
| 12 | inf | 연령대_40_상주인구_비율 |
| 13 | inf | 연령대_50_상주인구_비율 |
| 14 | inf | 연령대_60_이상_상주인구_비율 |
| 15 | 32.874815 | 총_가구_수 |
| 16 | 5.411849 | 집객시설_수 |
| 17 | 2.586857 | 월_평균_소득_금액 |
| 18 | 6.784772 | 지출_총금액 |
| 19 | 15.814548 | 유사_업종_점포_수 |
| 20 | 1.357144 | 개업_점포_수 |
| 21 | 1.417630 | 폐업_점포_수 |
| 22 | 1.574667 | 편의점_밀도 |

# CONTENT

**INTRODUCTION**

- TEAM MEMBERS
- BACKGROUND
- WEB SERVICE
- DEVELOPMENT
  ENVIRONMENT
- PROJECT PERIOD

**PROJECT**

- FLOW CHART
- WBS

**DATA**

- DATA COLLECTION
- DATA PREPROCESSING

**EDA**

- EDA
- CORRELATION
  ANALYSIS
- MULTICOLLINEARITY

**MODELING**

- MODEL
- MODEL TRAINING
- MODEL PERFORMANCE
  VALIDATION
- CONVENIENCE STORE
  SALES PREDICTION

**APP & DOCUMENT**

- STREAMLIT
- LIMITATIONS
  IMPROVEMENTS
- REFERENCES
- APPENDIX

# 5. MODELING

**TIME SERIES ANALYSIS**

- Analyzing **the change in data y over time** e.g.) predicting stock prices, sales, or temperatures

- Typically involves examining how it varies according to trends, cycles, seasonal components, and irregular or random elements. It's assumed that these components cause fluctuations in the data.

- Observing patterns of data variation → Segmenting them into trend, seasonal, and irregular components → Applying forecasting techniques like exponential modeling or ARIMA (AutoRegressive Integrated Moving Average) methods

**Trend:** A component that represents the overall upward or downward direction of observations.

**Cycle:** A component indicating changes that are periodic but not seasonal, characterized by longer cycles.

**Seasonal:** A component that represents factors explained by regular variations according to specific periods.

**Irregular:** A component representing errors that cannot be explained by specific patterns.

**Random:** A component representing random causes that are independent of regular movements over time.

# 5. MODELING

**ARIMA & SARIMA**

### ARIMA (AutoRegressive Integrated Moving Average)

- Suitable for modeling the trend and volatility of non-seasonal data.
- **Differencing** is used to address the non-stationarity of time series.

### SARIMA (Seasonal ARIMA)

- Extended version of ARIMA that can additionally model seasonal patterns.
- **Seasonal differencing** is used to handle the seasonality of time series.



## Model Performance

**ARIMA - MAE** 21,367,509,120.59 / **RMSE** 43,353,668,254.71

**SARIMA - MAE** 43,774,410,318.39 / **RMSE** 77,763,589,777.93

The ARIMA model shows a lower Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) compared to the SARIMA model.

**Since the error range is quite wide, considering other models is preferred.**

# 5. MODELING

**HOLT-WINTERS MODEL TIME SERIES PREDICTION**

- Statistical techniques used for analyzing and forecasting time series data
- Considers the three main components that may exist in time series data: **level**, **trend**, and **seasonality**.
- Predict future data points by accounting for all these components, useful for short-term forecasts of time series data



### Model Performance

**Mean Absolute Error (MAE) : 7,317,570,194.14**

**Root Mean Squared Error (RMSE) : 8,237,463,143.96**

**Since the error range is quite wide,
considering other models is preferred.**

- For convenience, assuming the first month of each quarter
- Fitting the Holt-Winters model/future predictions

# 5. MODELING

**FINAL MODEL SELECTION**

- After applying various machine learning algorithms such as time series, RandomForest, XGBoost, and LightGBM, the **LightGBM** model showed the best performance.

| Algorithm | Data Reprocessing | Feature Engineering | Cross validation | Hyperparameter Tuning | MAE | MSE | RMSE | R-Squared |
|---|---|---|---|---|---|---|---|---|
| LGBM | Categorical Variable : Dummy Encoding / Numeric Variable Standard Scaler Dependent Variable : Box-Cox Transformation | Proportion of Resident Population by Age Proportion of Working Population by Age Convenience Store Density | | num_leaves': [25, 30, 35], 'learning_rate': [0.12, 0.13, 0.14], 'n_estimators': [375, 400, 425] | 2,136,328,813,899,580 | 26930386.26 | 46220437.19 | **0.9885977983** |

# 5. MODELING

**DEPENDENT VARIABLE SCALING**

- Box-Cox transformation is a method of adjusting the distribution of data to control skewness.
- Commonly used to normalize data distributions in models where the data distribution is not normal.

### Log Transformation

- The dependent variable is left-skewed, so attempted log transformation to convert it to a normal distribution.
- However, after log transformation, it became right-skewed and the distribution issue was not resolved.



### Box-Cox Transformation

- Attempted Box-Cox transformation to address the skewness in the distribution.
- The dependent variable in both the train and test datasets transformed into a shape similar to a normal distribution.

# 5. MODELING

**FEATURE IMPORTANCE**

- Confirmed that the feature **[Hourly Floating Population]** which showed a high correlation in the correlation analysis is the most important feature.



Top 20 Most Important Features

# 5. MODELING

**Actual Values vs Predicted Values**

- When visualizing actual values versus predicted values, it is observed that they are distributed almost evenly around the y=x line.

  => This indicates that the model predicts the data well.

# CONTENT

## INTRODUCTION

- TEAM MEMBERS
- BACKGROUND
- WEB SERVICE
- DEVELOPMENT
  ENVIRONMENT
- PROJECT PERIOD

## PROJECT

- FLOW CHART
- WBS

## DATA

- DATA COLLECTION
- DATA PREPROCESSING

## EDA

- EDA
- CORRELATION
  ANALYSIS
- MULTICOLLINEARITY

## MODELING

- MODEL
- MODEL TRAINING
- MODEL PERFORMANCE
  VALIDATION
- CONVENIENCE STORE
  SALES PREDICTION

## APP & DOCUMENT

- STREAMLIT
- LIMITATIONS
  IMPROVEMENTS
- REFERENCES
- APPENDIX

# 6. DOCUMENTS

**SERVICE (STREAMLIT)**



**Menu Composition**

메뉴

🖥 메뉴를 선택
하세요:

⌂ 홈

▥ 강남구 편의점 분포
현황

📈 강남구 편의점 매출
현황

💳 매출 현황 순위

⚙ 매출 예측 모델링



**Streamlit - HOME**

강남구 편의점 매출 예측 🏪

이 앱은 편의점 예비
창업자들을 위한 **강남구**
지역의 **편의점** 시간대별
**매출 예측** 서비스입니다.

이 앱이 여러분들께
조금이나마 도움이 되기를
바랍니다.



**Streamlit - EDA**

강남구 편의점 매출 현황 📊

행정동 코드명 선택:

논현1동

선택된 행정동 코드명: 논현1동

논현1동의 시간대별 평균 매출

논현1동의 시간대별 평균 매출

강남구 편의점 예상매출 종합 🧠

예상매출 상권 추천    예상 순수익

시간대 희망매출 최대값(단위:백만)

시간대 희망매출 최소값(단위:백만)

# 6. DOCUMENTS

**LIMITATIONS / IMPROVEMENTS**

- Intended to implement a monthly net profit prediction model by calculating labor costs and franchise operating expenses, but unable to do so due to time constraints.

# 6. DOCUMENTS

REFERENCES

- **PAPERS**

  1) 김현철, 이승일, 2019, "서울시 골목상권 매출액에 영향을 미치는 요인에 관한 연구", 「서울도시연구」, 제 20권 제 1호
  2) 황규성, 2014, "편의점 입지선정시 매출에 영향을 미치는 요인분석"
  3) 김미성, 2020, "서울시 상권 데이터의 시각화에 기반한 매출액 예측"
  4) 이철환, 2012, "편의점의 상권 추정과 매출 예측에 관한 연구"
  5) 김동명, 2020, "시스템 다이내믹스를 활용한 편의점 특정 상품 매출 분석 및 예측"
  6) 이임동, 이찬호, 강상목, 2010, "편의점 매출에 영향을 미치는 입지요인에 대한 실증연구"

- **NEWS ARTICLES**

  1) 송금종, 카페 1950개·편의점 470개…강남구, 서울 최대 '슬세권', 쿠키뉴스, 2023.12.24
     https://www.kukinews.com/newsView/kuk202312140078

  2) 이지원, 편의점 본사가 제시한 '예상 매출액'이 과장이라면…, 더스쿠프, 2023.02.14
     https://www.thescoop.co.kr/news/articleView.html?idxno=56799

  3) 이진원, 홈플러스, 예비 편의점주에 예상매출 뻥튀기, 시민일보, 2017.11.05
     https://www.siminilbo.co.kr/news/articleView.html?idxno=537797

# 6. DOCUMENTS

**APPENDIX - SERVICE**

**Page Navigation Sidebar**

- Select a menu to navigate to the desired page



메뉴

메뉴를 선택
하세요:

홈

강남구 편의점 분포
현황

강남구 편의점 매출
현황

매출 현황 순위

매출 예측 모델링



강남구 편의점 매출 예측

이 앱은 편의점 예비
창업자들을 위한 **강남구**
지역의 **편의점** 시간대별
**매출 예측** 서비스입니다.

이 앱이 여러분들께
조금이나마 도움이 되기를
바랍니다.

# 6. DOCUMENTS

**APPENDIX - SERVICE**



## Page Navigation Sidebar

- Select a menu to navigate to the desired page

## Home Screen

- Purpose of the service

강남구 편의점 매출 예측

메뉴

📺 메뉴를 선택 하세요:

⌂ 홈

▥ 강남구 편의점 분포 현황

📈 강남구 편의점 매출 현황

🕹 매출 현황 순위

◉ 매출 예측 모델링

이 앱은 편의점 예비 창업자들을 위한 **강남구** 지역의 **편의점** 시간대별 **매출 예측** 서비스입니다.

이 앱이 여러분들께 조금이나마 도움이 되기를 바랍니다.

# 6. DOCUMENTS

**APPENDIX - SERVICE**

**Map Visualization**

- Checking the distribution status of the desired commercial area

# 6. DOCUMENTS

**APPENDIX - SERVICE**

## Map Visualization

- Checking the distribution status of the desired commercial area

## Detailed Information

- Average sales amount for the selected commercial area



메뉴

메뉴를 선택하세요:

⌂ 홈

🔲 강남구 편의점 분포 현황

📈 강남구 편의점 매출 현황

🏆 매출 현황 순위

◉ 매출 예측 모델링

강남구 편의점 분포 현황 🗺️
(2021년 1분기 ~ 2023년 3분기)

궁금한 상권을 선택하세요 👀

상권명 : 대청초등학교, 행정동 : 일원1동, 시간대_매출_금액_평균 : 502.41백만원

# 6. DOCUMENTS

**APPENDIX - SERVICE**

| Selectbox |
|---|
| - Select the desired administrative district |

# 6. DOCUMENTS

**APPENDIX - SERVICE**

## Selectbox

- Select the desired administrative district

## Bar Graph

- The hourly average sales for the administrative district selected by the user

강남구 편의점 매출 현황 📊

메뉴

메뉴를 선택 하세요:

🖥 홈

📖 강남구 편의점 분포 현황

📈 강남구 편의점 매출 현황

🎯 매출 현황 순위

⚙ 매출 예측 모델링

행정동 코드명 선택:

논현1동

선택된 행정동 코드명: 논현1동

논현1동의 시간대별 평균 매출

논현1동의 시간대별 평균 매출

893011669.23
703400202.68
695708760.75
522646817.85
459139463.06
478443683.04

평균매출금액

8억
6억
4억
2억
0억

00~06  06~11  11~14  14~17  17~21  21~24

시간대

# 6. DOCUMENTS

**APPENDIX - SERVICE**

**Selectbox**

- Select the commercial area corresponding to the administrative district chosen above

메뉴

🖵 메뉴를 선택 하세요:

⌂ 홈

▥ 강남구 편의점 분포 현황

✓ **강남구 편의점 매출 현황**

🖩 매출 현황 순위

◉ 매출 예측 모델링

행정동 코드명 선택:

논현1동 ⌄

선택된 행정동 코드명: 논현1동

논현1동의 시간대별 평균 매출 ⌄

논현1동에 대한 상권 코드명 선택:

논현초등학교 ⌄

논현초등학교 상권의 시간대별 평균 매출 ⌃


논현초등학교 상권의 시간대별 평균 매출

# 6. DOCUMENTS

**APPENDIX - SERVICE**

행정동 코드명 선택:

논현1동

선택된 행정동 코드명: 논현1동

논현1동의 시간대별 평균 매출

**메뉴**

메뉴를 선택 하세요:

홈

강남구 편의점 분포 현황

강남구 편의점 매출 현황

매출 현황 순위

매출 예측 모델링

논현1동에 대한 상권 코드명 선택:

논현초등학교

**Selectbox**

- Select the commercial area corresponding to the administrative district chosen above

논현초등학교 상권의 시간대별 평균 매출

논현초등학교 상권의 시간대별 평균 매출

539375234.82

483066288.73

5억

378710231.27

4억

평균매출금액

3억

267416751.82

240345718.55

260302267.55

2억

1억

0억

00~06   06~11   11~14   14~17   17~21   21~24

시간대

**Bar Graph**

- The hourly average sales for the commercial area selected by the user

# 6. DOCUMENTS

**APPENDIX - SERVICE**



**Radio Button**

- Select the desired time slot

# 6. DOCUMENTS

**APPENDIX - SERVICE**

| Radio Button |
|---|
| - Select the desired time slot |

| Bar Graph |
|---|
| - The sales of the top 5 commercial areas that recorded the highest sales during the time slot selected by the user |

메뉴

📺 메뉴를 선택
하세요:

🏠 홈

🎛 강남구 편의점 분포
현황

📈 강남구 편의점 매출
현황

🔲 매출 현황 순위

⚙ 매출 예측 모델링

**매출 현황 순위 💰**

시간대를 선택하세요:
- 🔴 00:00 ~ 06:00
- ⚪ 06:00 ~ 11:00
- ⚪ 11:00 ~ 14:00
- ⚪ 14:00 ~ 17:00
- ⚪ 17:00 ~ 21:00
- ⚪ 21:00 ~ 24:00



00:00 ~ 06:00 시간대 매출이 가장 높은 상권 TOP5

1410271160.82 · 891056281.09 · 876115064.82 · 633725267.18 · 498861772.91

# 6. DOCUMENTS

**APPENDIX - SERVICE**

# 6. DOCUMENTS

**APPENDIX - SERVICE**

| Radio Button |
|---|
| - Select the desired type |

| Selectbox |
|---|
| - Select the desired administrative district, commercial area, and quarter |

# 6. DOCUMENTS

**APPENDIX - SERVICE**

| Radio Button |
| --- |
| - Select the desired type |

| Selectbox |
| --- |
| - Select the desired administrative district, commercial area, and quarter |

| Bar Graph |
| --- |
| - The hourly projected sales based on user-selected conditions |

# 6. DOCUMENTS

**APPENDIX - SERVICE**



| Tab |
| --- |
| - Move to the desired tab |

# 6. DOCUMENTS

**APPENDIX - SERVICE**

| Tab |
| --- |
| - Move to the desired tab |

| Selectbox |
| --- |
| - Select the desired time slot |

# 6. DOCUMENTS

**APPENDIX - SERVICE**



| Tab |
| --- |
| - Move to the desired tab |

| Selectbox |
| --- |
| - Select the desired time slot |

| Detailed Information |
| --- |
| - The projected sales amount for the time slot selected by the user |

# 6. DOCUMENTS

**APPENDIX - SERVICE**

| Tab |
| --- |
| - Check the projected sales for each commercial area by time slot, year, and month individually |

# 6. DOCUMENTS

**APPENDIX - SERVICE**

## Tab

- Check the projected sales for each commercial area by time slot, year, and month individually

## Detailed Information

- Projected sales by time slot, year, and month for each commercial area can be visualized on a map

# 6. DOCUMENTS

**APPENDIX - SERVICE**

| Slider |
| :---: |
| - The user inputs the desired minimum and maximum values for sales |

# 6. DOCUMENTS

**APPENDIX - SERVICE**

## Slider

- The user inputs the desired minimum and maximum values for sales

## Detailed Information

- Information about commercial areas that meet the user's input values on the map

강남구 편의점 예상매출 종합 🧠

예상 매출 상권 추천    예상 순 수익

월 희망매출 최대값(단위:백만)
3000
0                                                    3000

월 희망매출 최소값(단위:백만)
801
0                                                    3000

상권명: 선릉역
행정동: 대치4동
예상 월 매출 금액(단위:백만): 1513.83

메뉴

📺 메뉴를 선택하세요:

🏠 홈

🏬 강남구 편의점 분포 현황

📈 강남구 편의점 매출 현황

🏷 매출 현황 순위

⚙ 매출 예측 모델링

유형을 선택하세요
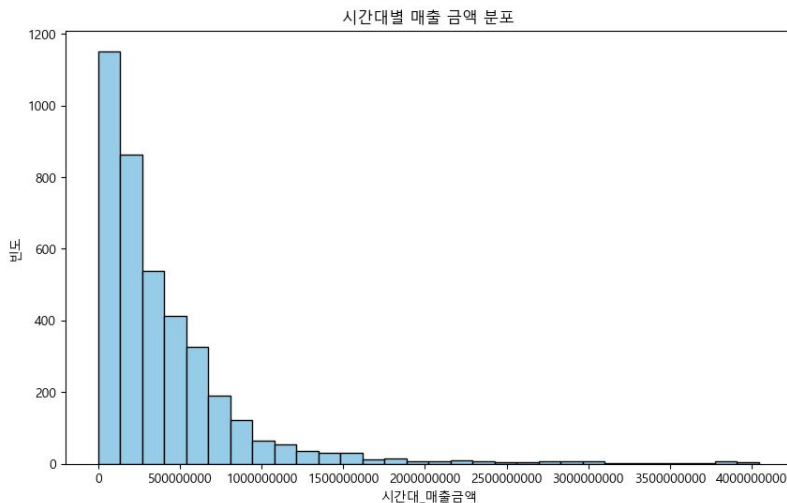◯ 선택 상권 예상 매출
◯ 예상 매출 종합
◉ 예상 매출 상권 추천

# 6. DOCUMENTS

APPENDIX - EDA

**Continuous Dependent Variable**



시간대별 매출 금액 분포

| Final Columns | |
|---|---|
| - Sales Amount by Time Period | - Monthly Average Income Amount |
| - Base Year | |
| - Base Quarter | - Total Expenditure Amount |
| - Commercial District Division Name | - Number of Stores in Similar Industries |
| - Commercial District Name | - Number of Newly Opened Stores |
| - Administrative Neighborhood Name | - Number of Closed Stores |
| - Time Period | |
| - Floating Population Count by Time Period | |
| - Total Working Population | |
| - Total Resident Population | |
| - Total Number of Households | |
| - Number of Attraction Facilities | |

# 6. DOCUMENTS

**APPENDIX - EDA**



**Continuous Independent Variables**

**Final Columns**

- Sales Amount by Time Period
- Base Year
- Base Quarter
- Commercial District Division Name
- Commercial District Name
- Administrative Neighborhood Name
- Time Period
- Floating Population Count by Time Period
- Total Working Population
- Total Resident Population
- Total Number of Households
- Number of Attraction Facilities

- Monthly Average Income Amount
- Total Expenditure Amount
- Number of Stores in Similar Industries
- Number of Newly Opened Stores
- Number of Closed Stores

# 6. DOCUMENTS

APPENDIX - EDA

Continuous Independent Variables
But, Similar to Categorical Data Distribution
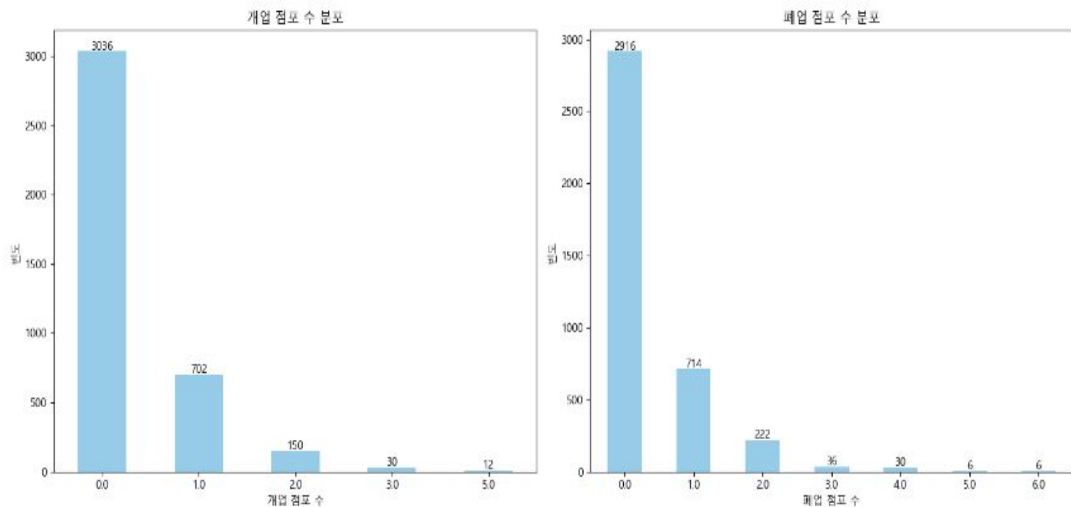


개업 점포 수 분포



폐업 점포 수 분포

### Final Columns

| | |
|---|---|
| - Sales Amount by Time Period | - Monthly Average Income Amount |
| - Base Year | - Total Expenditure Amount |
| - Base Quarter | - Number of Stores in Similar Industries |
| - Commercial District Division Name | - Number of Newly Opened Stores |
| - Commercial District Name | - Number of Closed Stores |
| - Administrative Neighborhood Name | |
| - Time Period | |
| - Floating Population Count by Time Period | |
| - Total Working Population | |
| - Total Resident Population | |
| - Total Number of Households | |
| - Number of Attraction Facilities | |

# 6. DOCUMENTS

**APPENDIX - EDA**

**Categorical Independent Variables**



**Final Columns**

- Sales Amount by Time Period
- Base Year
- Base Quarter
- Commercial District Division Name
- Commercial District Name
- Administrative Neighborhood Name
- Time Period
- Floating Population Count by Time Period
- Total Working Population
- Total Resident Population
- Total Number of Households
- Number of Attraction Facilities
- Monthly Average Income Amount
- Total Expenditure Amount
- Number of Stores in Similar Industries
- Number of Newly Opened Stores
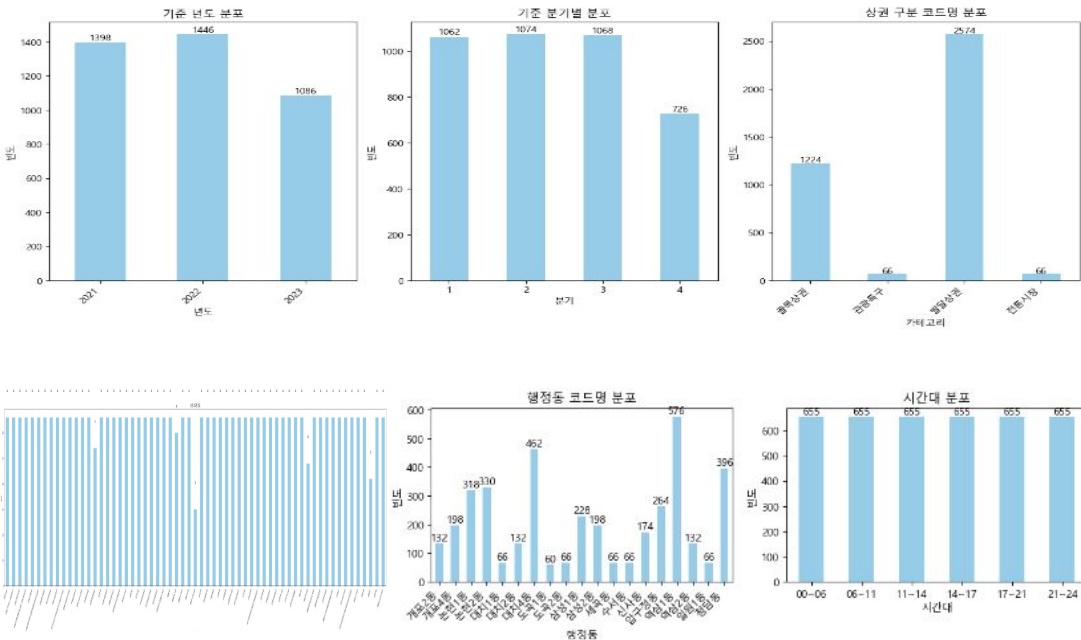- Number of Closed Stores

# 6. DOCUMENTS

**APPENDIX - MODELING**

| Algorithm | Data Reprocessing | Feature Engineering | Cross validation | Hyperparameter Tuning | MAE | MSE | RMSE | R-Squared |
|---|---|---|---|---|---|---|---|---|
| RandomForest | StandardScaler | Employment Population Ratio by Age Group | | | | | 67898424.45 | **0.9648454212** |
| RandomForest | StandardScaler | Ratio of Household Count to Population Count | | | | | 66123873.23 | **0.9754612381** |
| RandomForest | OneHotEncoder | Employment Population Ratio by Age Group | | | | | 66731521.15 | **0.9648153123** |
| RandomForest | OneHotEncoder | Ratio of Household Count to Population Count | | | | | 66197315.94 | **0.9618132176** |

# 6. DOCUMENTS

**APPENDIX - MODELING**

| Algorithm | Data Reprocessing | Feature Engineering | Cross validation | Hyperparameter Tuning | MAE | MSE | RMSE | R-Squared |
|---|---|---|---|---|---|---|---|---|
| RandomForest | StandardScaler: Numeric Data OneHotEncoder: Categorical Data | Ratio of Household Count to Population Count | KFold | | | | 69963979.84 | **0.9798526552** |
| RandomForest | StandardScaler: Numeric Data OneHotEncoder: Categorical Data | Ratio of Household Count to Population Count | KFold | {'n_estimators': 500, 'min_samples_split': 2, 'min_samples_leaf: 1, 'max_depth': 50, 'bootstrap': True} | | | 69628037.42 | **0.9800303186** |
| RandomForest + GradientBoosting | StandardScaler: Numeric Data OneHotEncoder: Categorical Data | StandardScaler: Numerical Data OneHotEncoder: Categorical Data | | | | | | **0.9773915991** |

# 6. DOCUMENTS

**APPENDIX - MODELING**

| Algorithm | Data Reprocessing | Feature Engineering | Cross validation | Hyperparameter Tuning | MAE | MSE | RMSE | R-Squared |
|-----------|-------------------|---------------------|------------------|----------------------|-----|-----|------|-----------|
| XGBoost | StandardScaler: Numeric Data OneHotEncoder: Categorical Data | Ratio of Household Count to Population Count, Ratio of Floating Population on Weekdays to Weekends | KFold | | | | 60577427.09 | **0.9849553596** |
| XGBoost | StandardScaler: Numeric Data OneHotEncoder: Categorical Data | Ratio of Household Count to Population Count, Ratio of Floating Population on Weekdays to Weekends | KFold | {'subsample': 0.6, 'n_estimators': 500, 'min_child_weight': 1, 'max_depth': 7, 'learning_rate': 0.05, 'colsample_bytree': 0.8} | | | 57367659.29 | **0.9864766996** |
| XGBoost | Removing 'Hourly Sales Amount', 'Average Weekday Floating Population Count', 'Average Weekend Floating Population Count', 'Total Resident Population Count', 'Commercial Area Code', 'Administrative District Code Name' | Ratio of Household Count to Population Count | | | | | 75272170.99 | **0.9812817025** |

# 6. DOCUMENTS

**APPENDIX - MODELING**

| Algorithm | Data Reprocessing | Feature Engineering | Cross validation | Hyperparameter Tuning | MAE | MSE | RMSE | R-Squared |
|---|---|---|---|---|---|---|---|---|
| LGBM | Categorical Variable One-Hot Encoding / Numeric Variable Standard Scaling Dependent Variable Box-Cox Transformation | Proportion of Resident Population by Age Proportion of Working Population by Age Convenience Store Density | | num_leaves=31, learning_rate=0.1, n_estimators=100 | 3,518,915,356, 168,770 | 33206245.03 | 59617671 | **0.9812185361** |
| LGBM | Categorical Variable One-Hot Encoding / Numeric Variable Standard Scaling Dependent Variable Box-Cox Transformation | Proportion of Resident Population by Age Proportion of Working Population by Age Convenience Store Density | | GridSearchCV 'num_leaves': [15, 31, 50], 'learning_rate': [0.05, 0.1, 0.2], 'n_estimators': [50, 100, 200] | 2,283,010,610, 939,760 | 27001065.95 | 47780860 | **0.9878149154** |
| LGBM | Categorical Variable One-Hot Encoding / Numeric Variable Standard Scaling Dependent Variable Box-Cox Transformation | Proportion of Resident Population by Age Proportion of Working Population by Age Convenience Store Density | | num_leaves': [15, 31, 50], 'learning_rate': [0.1, 0.15, 0.2], 'n_estimators': [200, 300, 400] | 2,289,286,910, 788,070 | 27355461.16 | 47846493 | **0.987781417** |

# THANK YOU