

# Urban Polarization – Data Analysis



# Contents

- I. Project Overview
- I. Team member introduction and roles
- I. Development process
- I. conclusion

# What is urban polarization? I asked GPT!



“Occurs within or between cities  
“It refers to the phenomenon of deepening

**Urban Polarization** refers to the process by which economic, social, spatial, or demographic differences within a city become increasingly divided into contrasting extremes, rather than remaining evenly distributed across the urban population or territory.

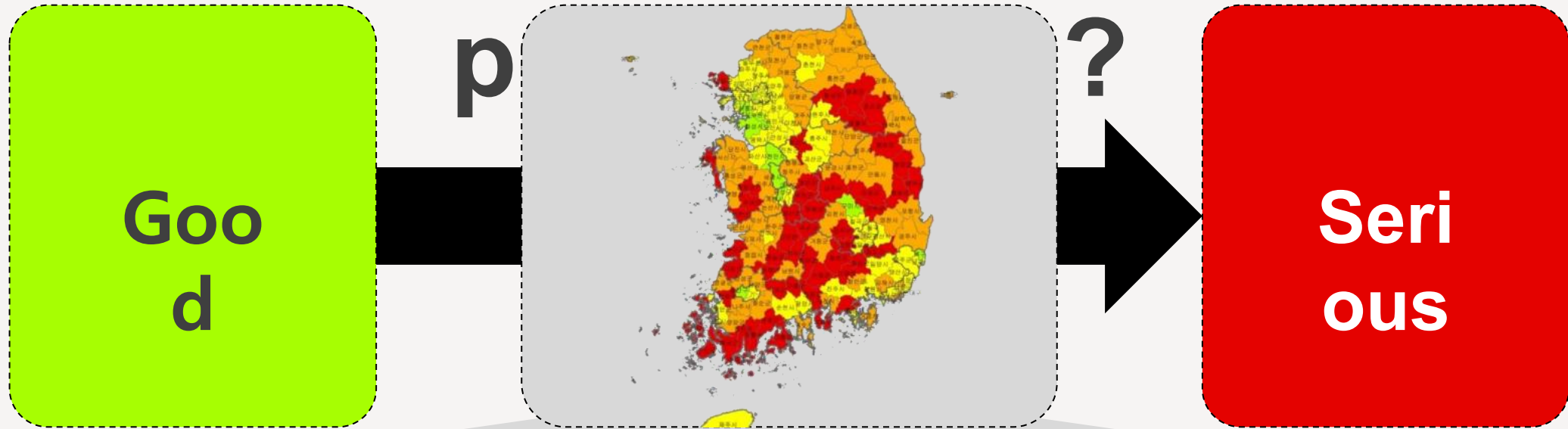
In practice, urban polarization typically manifests as:

- **Economic polarization:** A growing gap between high-income and low-income residents, often accompanied by the expansion of high-wage professional jobs alongside low-wage service work, with a shrinking middle class.
- **Spatial polarization:** The physical separation of social groups within the city, such as affluent neighborhoods coexisting with areas of concentrated poverty, informal settlements, or marginalized districts.
- **Social and demographic polarization:** Segregation based on factors such as education, ethnicity, immigration status, or access to public services.
- **Access polarization:** Unequal access to housing, transportation, healthcare, education, and employment opportunities across different urban areas.

Urban polarization is commonly linked to processes such as globalization, deindustrialization, gentrification, housing market dynamics, and uneven urban policy interventions. Over time, it can intensify social inequality, reduce social cohesion, and create persistent disparities in life chances within the same city.



## Why



“The society is seriously aware of the problems of regional extinction and local extinction. Accordingly, we hope to help resolve

- Analysis subjects: Cities in the Chungcheong region (Sejong, Daejeon, North Chungcheong, South Chungcheong)

## Population influx



## Reason for selection



“It’s not a metropolitan area,  
but  
with

cities in danger of extinction

## Some extinction threats



it is considered appropriate to

# How to analyze urban polarization?

As expected... I asked GPT!



## economic

gap

“Rich areas and poor  
areas

“Regional income  
disparities, employment

## social gap

“Differences in  
access to social  
resources such as  
education, health,

## spatial

gap

“Imbalance in the  
distribution of  
residential,  
commercial, and

- statistical data


Proportion of  
the working-age  
population

crime rate

Apartment  
sales  
Actual

Average  
monthly income  
distribution

com  
paris  
on  
yze



## concl usion

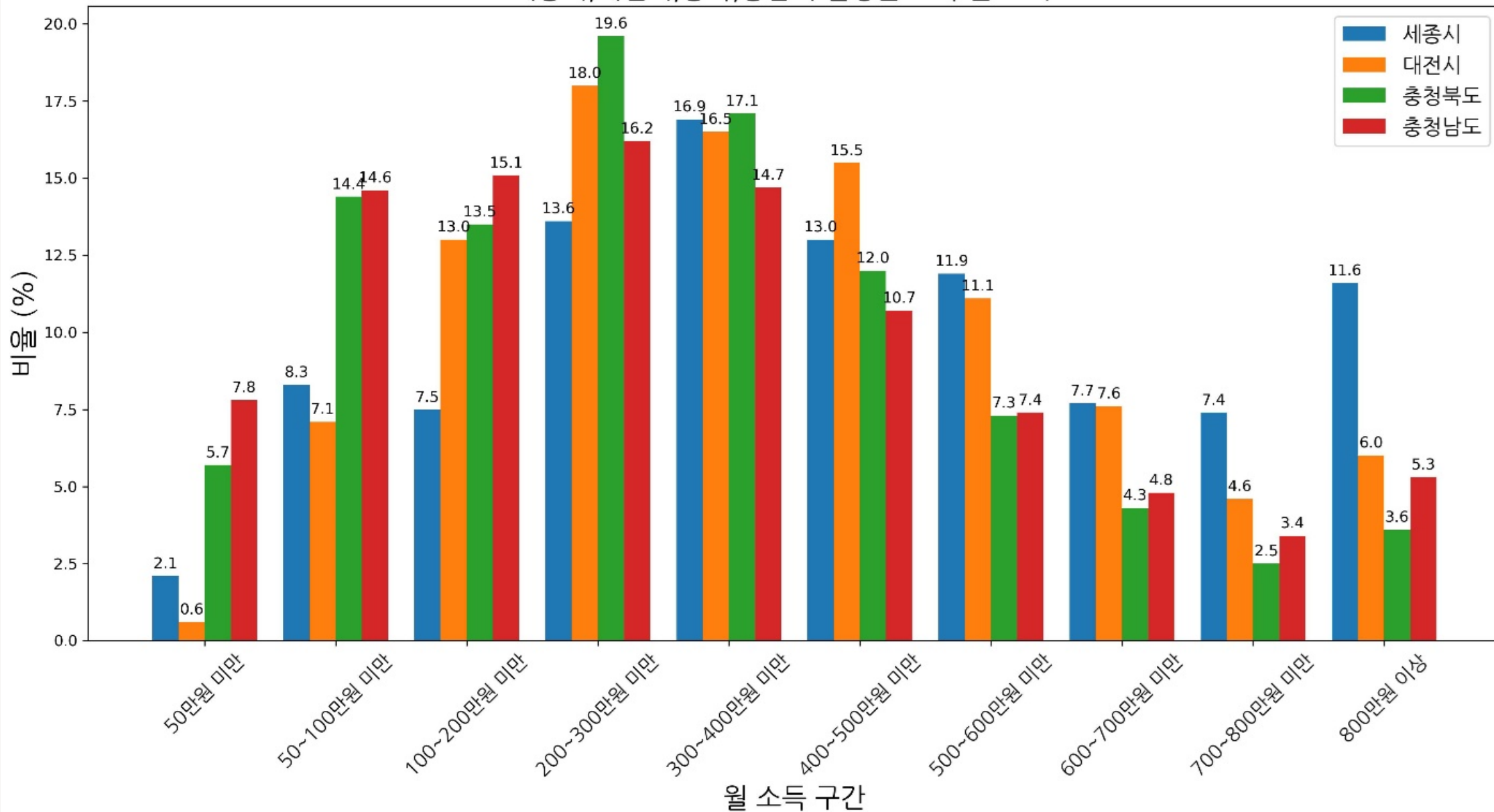
the working-age population ratio,  
apartment transaction price index, and  
monthly average income distribution,  
we can see that polarization is clearly  
evident.

Among them, we will select the apartment

transaction price index as a dependent



세종시,대전시,충북,충남의 월평균 소득 분포 비교



## Why did you choose apartment

prices as a key variable?  
Apartment prices are considered to be a variable that

clearly shows the relative wealth and poverty (degree of backwardness) of a city .

who make up the majority of the region's average monthly income , choose apartments that offer convenience or

luxury apartments as their primary residences, while it is

difficult for low-income people to enter apartments .

Variable  
selection  
method  
Backward

elimination method  
(Backward  
Elimination)

Initially there were 32  
variables

I trained using it, but  
Overfitting occurs

① Dimension  
reduction to solve  
(Variable removal,

Independent  
variables (final  
selection)

Regional  
population  
movement data  
(Total inflow, total

PIR  
index  
(transfer)

Comprehens  
ive real  
estate tax

Apartment  
sales volume

Apartment sales  
supply and  
demand trends

Presidential  
approval  
rating

Home mortgage  
interest rates

consumer  
price index

depend  
ent  
variabl  
e

Apartments by  
region

Real transaction  
index data

## Excluded variables and reasons

Number of single-person

- ◆ High correlation was confirmed, but quick reflection was not possible due to lack of monthly data.

Dominance in the number of National

- ◆ I tried inserting it using One-Hot Encoding, but the prediction rate decreased.

Housing Affordability Index

- ◆ Only quarterly data, not monthly data, is available, and quick reflection is not possible.
- ◆ Overlapping variables used with the PIR index

Base interest rate

- ◆ Since they show a strong positive correlation of 0.86, to resolve the multicollinearity problem

Adoption of mortgage loan interest rates with a high correlation with

LIR index

- ◆ The two assess the \*affordability of the housing market from different perspectives, The PIR index, which has a very high correlation of 0.96 and is highly correlated with the real transaction index, was adopted.

\* Generally, the extent to which an individual or household can afford the costs required to purchase or rent a home.

Data Analytics Professional Team

Seong



Data  
crawling  
Data

preprocessi

SHIN



DB design  
and  
construction

Data

Park



Planning and  
overall  
management

Building MI

Lee



Backend  
construction  
Data

visualization

Oh



Backend  
construction  
Data

visualization

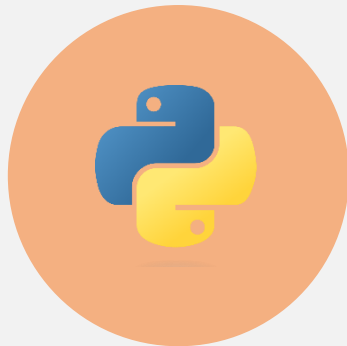
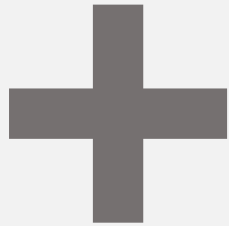




## Data collectio

n

### Case ① Crawling (automatic collection)



API  
(data URL)

crawling  
(URL-based  
collection)

### Case ② Research (direct collection)



API  
(Unavaila  
ble)



Access the  
webpage  
(Direct  
collection)



## Data collection

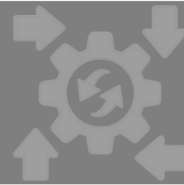
HF 한국주택금융공사  
주택금융연구원

Gallup 한국갤럽

Troubleshooting some variable crawl failures.

Failed to crawl data from the Korea Housing  
Finance Research Institute and Gallup Korea.

We plan to design and develop a crawling



## Data preprocess ing

### Method 1: Data frame defect based on the 'time' column

Point in time (reference column)	population movement data	PIR index
January 2013	100000	100
February 2013	200000	110
March 2013	300000	120
...	...	...
December 2023	1000000	140

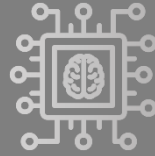
### ② Scaling

#### Min-Max scaling

$$x_{scaled} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

#### Robust scaling

$$X_{scaled} = \frac{X_i - \min(X)}{\max(X) - \min(X)}$$



## Model building

(Existing  
model)

Mobile  
data

8

variables

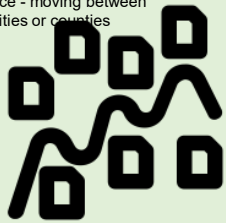
Dimens  
ion  
reducti  
on

4

variables

Total transfer  
Total transfer  
as movement  
Inter-city transfer  
interprovincial transfer  
Moving within the city/county  
Moving within a city or  
province - moving between  
cities or counties  
Moving within a city or  
province - moving between  
cities or counties

Total transfer  
Total transfer  
Inter-city transfer  
Intercity transfer



Overfitting  
occurs!

(result)



Bad  
!

**Recognize the need for additional independent variables!**

Prediction rate for the national apartment sales transaction index

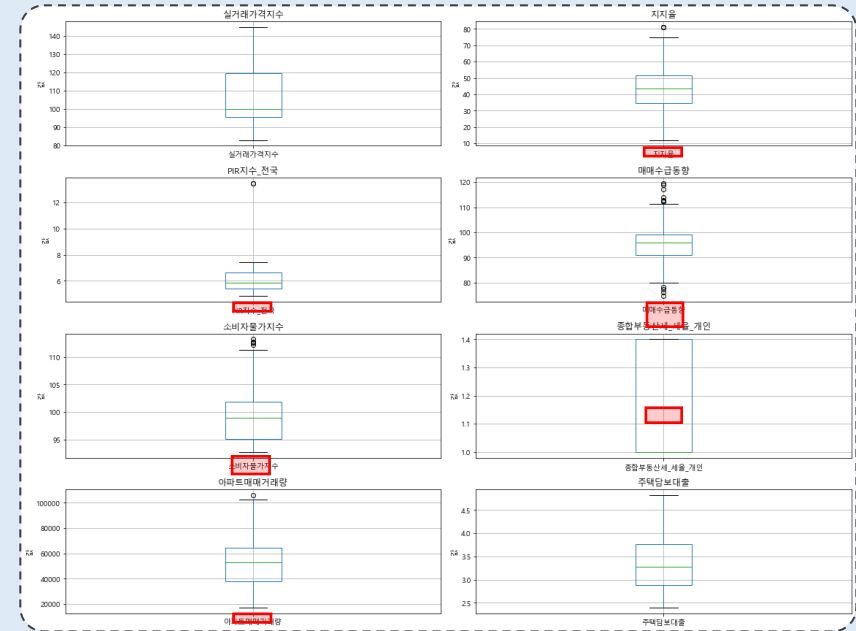
Secure and **mix with regional data to improve performance**

## ① Min-Max

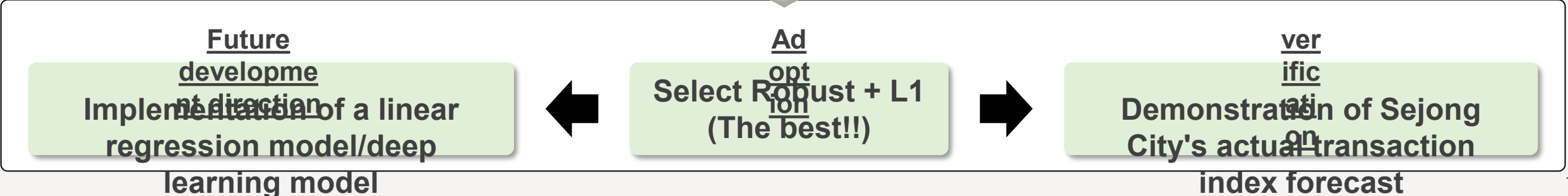
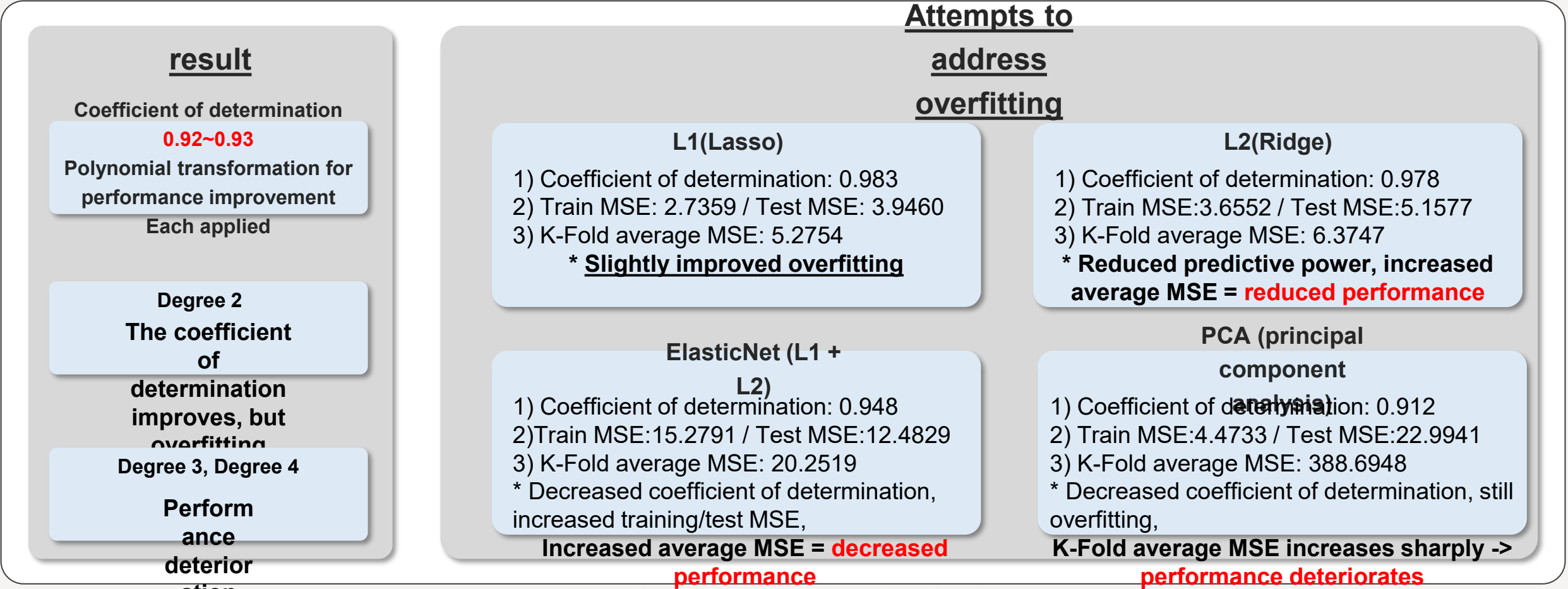
While maintaining data  
distribution  
Select to

## ② Robust

Boxbeard drawing  
Check variables where

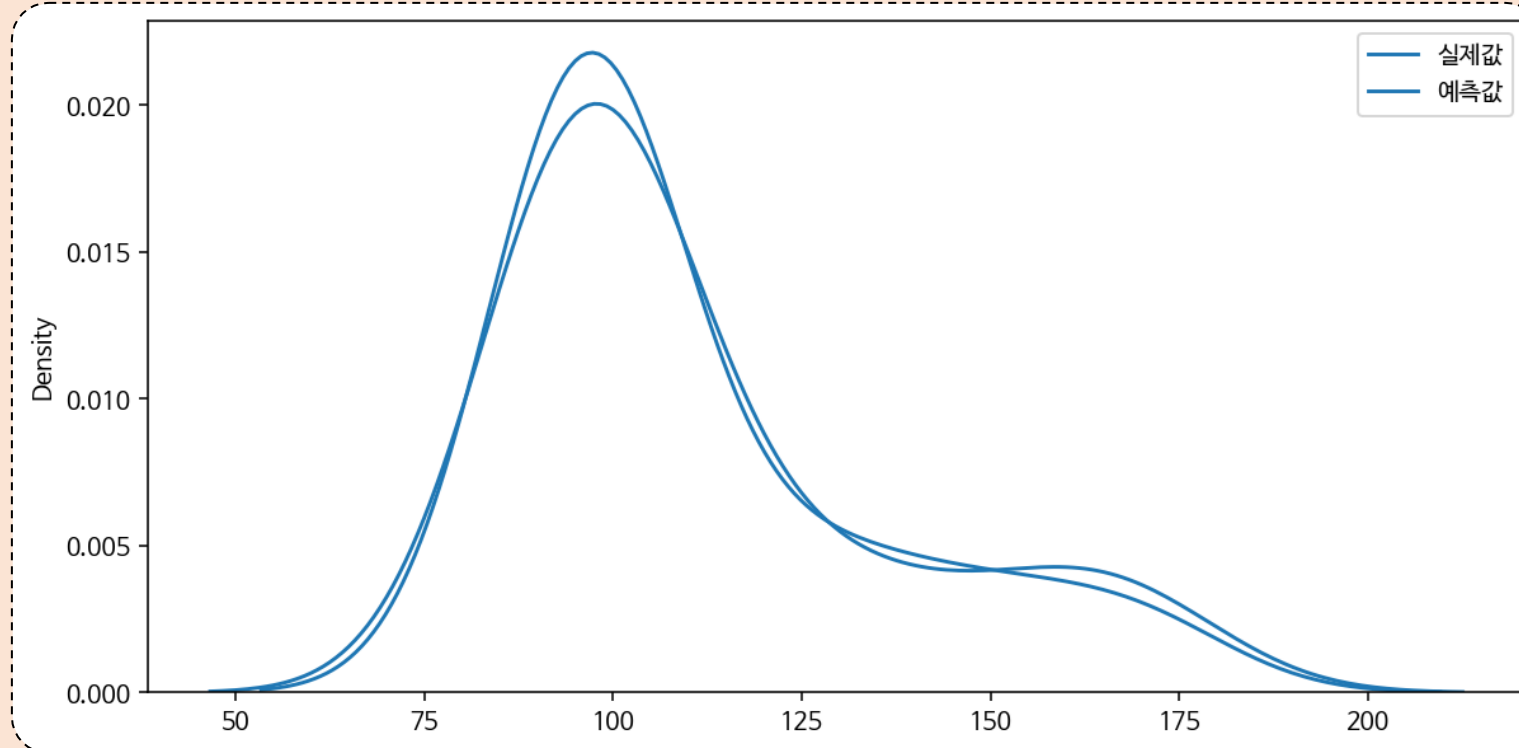






## Future Direction (still needs

We plan to improve performance by testing





## Data Visualization & Service Building



### Data Visualization & Service Building

- ① Net population migration: Implementation of a nationwide net population migration bar chart and map based on the combined data from 2013 to 2023.
- ② Real Transaction Price Index and Unemployment Rate: Insert the 2023 Daejeon/Sejong/Chungbuk/Chungnam Real Transaction Price Index and Unemployment Rate BAR CHART (image)
- ③ Crime Rate and Working-Age Population: Insert the 2023 Daejeon/Sejong/North Chungcheong/South Chungcheong crime rate and working-age population BAR CHART (image)



matplotlib



Streamlit



## Data Visualization & Service Building



Network URL: <http://192.168.71.220:8501>

#### 1. General

- 1) Map implementation: Daejeon, Sejong, South Chungcheong, and North Chungcheong regions are marked with colors.
- 2) Implementation of the Apartment Transaction Price Index LINE CHART: 4 regions: Daejeon, Sejong, Chungnam, and Chungbuk
- 3) Monthly average income bracket ratio implemented as a BAR CHART: 4 regions: Daejeon, Sejong, South Chungcheong, and North Chungcheong

#### 2. Real Transaction Index Prediction Model (Demo Version)

- 1) Select “Region,” “Start Date,” and “End Date” from the left menu.
- 2) Predicting future apartment transaction indices using linear regression machine learning algorithms.

#### 3. Real Estate Listing Recommendation Model (Demo Version)



## DB design and construction



All data loaded into DB

가족(Family)		
PK	fm_single	-1~4인가구
PK	fm_two	-4인가구
PK	fm_three	-4인가구
PK	fm_four	이상
PK	fm_fourmore	

실거래지수(Transaction_Price_index)		
PK	tp_Transaction_Price_index	실거래지수
FK	as_saving	날짜
FK	iv_date	

독립변수(Independent_variable)		
FK	Transaction_Price_index	
PK	President_support_rate	
PK	iv_date	
	PIR_national_index	
	Consumer_price_index	
	realestate_tax_for_individual	
	total_numberof_transaction	
	mortgage_loan	
	Supply_demands_trend	
	inflow_total	
	net_movement	
	outflow_total	

- 실거래지수
- 대통령 지지율
- 날짜(2013~2023 년월)
- PIR 지수
- 소비자 물가지수
- 종합부동산세(개인)
- 총거래량
- 주택담보대출
- 매매수급동향
- (총)전입
- (총)전출
- 순이동

- 나이
- 자산
- 대출
- 가족 구성원 수

나이(Age)		
PK	a_10s	INTEGER(2)
PK	a_20s	
PK	a_30s	
PK	a_40s	
PK	a_50s	
PK	a_60s	
PK	a_70s	
PK	a_80s	
PK	a_90s	
PK	a_100s	INTEGER(3)

Customer(고객)		
PK	cs_age	CHAR(3)
FK	as_savings	CHAR(10)
FK	as_loan	CHAR(10)
FK	fm_single	CHAR(1)
FK	fm_two	
FK	fm_three	
FK	fm_four	
FK	fm_fourmore	
FK	as_investment	

인프라_접근성(infrastructure-accessibility)		
PK	ia_transport	INTEGER
PK	ia_safety	
PK	ia_health	
PK	ia_education	
PK	ia_leisure	
FK	cs_age	CHAR(3)

- 교통
- 안전강화
- 교육
- 여가

자산(Assets)		
PK	as_savings	CHAR(10)
PK	as_loan	
PK	as_investment	
FK	cs_age	INTEGER

- 예금
- 대출
- 현금화 가능한 투자금 (전세금/주식/펀드/코인)





## DB design and construction



### Extracting data from a DB

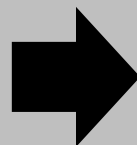
#### ① Writing a sample SQL query

# 머신러닝 데이터 선정

# 예시 : 2014.01 이전 데이터 + 인프라 수준 1등급 + 용산구

```
SELECT * FROM parameter where Pm_Date < 2014.01;
select inflow_total, outflow_total, net_movement from seoul where Io_Date < 2014.01;
select Subway, Primary_School, Middle_School, High_School, General_Hospital, Supermarket, Park
from fancy where district = 'yongsan-gu';

SELECT
p.*,
s.inflow_total, s.outflow_total, s.net_movement,
f.Subway, f.Primary_School, f.Middle_School, f.High_School, f.General_Hospital, f.Supermarket, f.Park
FROM parameter p
JOIN seoul s ON p.Pm_Date = s.Io_Date
JOIN fancy f ON f.district = 'yongsan-gu'
WHERE p.Pm_Date < '2014-01-01' AND s.Io_Date < '2014-01-01';
```



#### ② Data extraction

outflow_total	net_movement	Subway	Primary_School	Middle_School	High_School	General_Hospital
151771	-14152	981	1261	140	159	1200
131098	-11357	981	1261	140	159	1200
142141	-8773	981	1261	140	159	1200
109528	-6519	981	1261	140	159	1200
126101	-6930	981	1261	140	159	1200

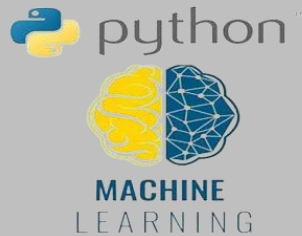


## DB design and construction

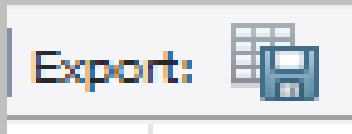


### Export and Model Training

❌ plan to write a module to extract data from the DB in the future.



#### ③ Export to CSV



파일 이름(N): train\_data.csv

파일 형식(T): CSV (\*.csv)

폴더 숨기기

저장(S) 취소

#### ④ Data for model

```
Pm_Date,Transaction_Price_index,President_support_rate,PIR_national_index,Consumer_price_index,realstat
2013.12,85.7,54,5.1,93.229,1.0,67883,3.74,"94.53",137619,151771,-14152,981,1261,140,159,1200,1400,1600
2013.11,85.7,54,5.04,93.116,1.0,61844,3.77,"95.61",119741,131098,-11357,981,1261,140,159,1200,1400,1600
2013.10,85.7,54,5.09,93.134,1.0,65871,3.81,"95.66",133368,142141,-8773,981,1261,140,159,1200,1400,1600
2013.09,85.2,60,5.08,93.419,1.0,39801,3.82,"92.23",103009,109528,-6519,981,1261,140,159,1200,1400,1600
2013.08,84.6,60,5.09,93.238,1.0,30794,3.8,"89.56",119171,126101,-6930,981,1261,140,159,1200,1400,1600
2013.07,84.3,60,5.15,92.909,1.0,25079,3.77,"90.41",118420,127308,-8888,981,1261,140,159,1200,1400,1600
2013.06,83.7,51,5.52,92.71,1.0,94647,3.73,"92.1",115773,124523,-8750,981,1261,140,159,1200,1400,1600
2013.05,84.0,51,5.62,92.823,1.0,64538,3.77,"93.34",130619,138830,-8211,981,1261,140,159,1200,1400,1600
2013.04,83.8,51,5.65,92.823,1.0,55442,3.86,"92.34",134581,144586,-10005,981,1261,140,159,1200,1400,1600
2013.03,83.5,42,5.53,92.952,1.0,47375,3.97,"88.63",138507,146220,-7713,981,1261,140,159,1200,1400,1600
2013.02,83.1,42,5.49,93.038,1.0,34089,4.06,"87.21",150805,155528,-4723,981,1261,140,159,1200,1400,1600
2013.01,83.0,42,5.35,92.728,1.0,16968,4.17,"86.1",118477,123006,-4529,981,1261,140,159,1200,1400,1600
```

## Why You Should Use a DBMS to Build a Machine Learning

## Model (1)

## SCHEMAS

Filter objects

haksa

▼ infrastructure

▼ Tables

▶ busan

▶ cheap

▶ chungnam

▶ chungbuk

▶ daegu

▶ daejeon

▶ fancy

▶ jeonnam

▶ jeonbuk

▶ jeju

▶ incheon

▶ gyeongbuk

▶ gwangju

▶ gangwondo

▶ parameter

▶ sejoong

DBMS

## Seoul City's High-Rise Apartment

## Infrastructure Data – fancy

district	Subway	Primary_School	Middle_School	High_School	General_Hospital
yongsan-gu	302	769	772	741	1900
eunpyeong-gu	363	462	763	893	2600
jonglo-gu	499	756	1200	1200	1000
jung-gu	310	1417	286	224	1800
junglang-gu	845	601	856	491	860
seongdong-gu	288	428	840	336	1900

## Seoul's low-cost apartments -

## infrastructure data – cheap

district	Subway	Primary_School	Middle_School	High_School	General_Hospital
yongsan-gu	981	1261	140	159	1200
eunpyeong-gu	413	367	344	929	3000
jonglo-gu	588	1000	382	365	430
jung-gu	514	524	139	183	1500
junglang-gu	575	733	983	857	342
seongdong-gu	528	706	619	572	2300

## Why You Should Use a DBMS to Build a Machine Learning Model (2)

### SCHEMAS

Filter objects

- haksa
- ▼ infrastructure
  - ▼ Tables
    - busan
    - cheap
    - chungnam
    - chungbuk
    - daegu
    - daejeon
    - fancy
    - jeonnam
    - jeonbuk
    - jeju
    - incheon
    - gyeongbuk
    - gwangju
    - gangwondo
    - parameter
    - sejoong

DBMS

### Variables with a high contribution to the real transaction index – parameter

Pm_Date	Transaction_Price_in	President_support_ra	PIR_national_index	Consumer_price_index
2013.01	83.0	42	5.35	92.728
2013.02	83.1	42	5.49	93.038
2013.03	83.5	42	5.53	92.952
2013.04	83.8	51	5.65	92.823
2013.05	84.0	51	5.62	92.823
2013.06	83.7	51	5.52	92.71

### Population Migration

Io_Date	inflow_total	outflow_total	net_movement
2013.01	39348	40048	-700
2013.02	42857	42807	50
2013.03	39546	41004	-1458
2013.04	40603	41548	-945

## Data extraction format (planning to use SQL)

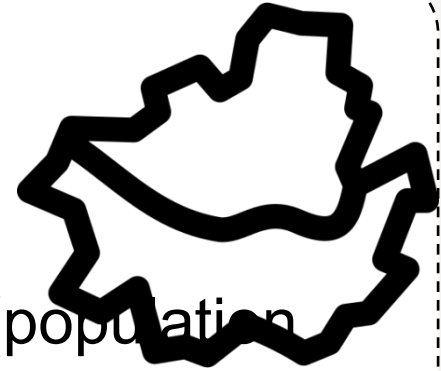
Example of use)

1) 2013~2014 data + infrastructure level 1 (Seocho-gu) + Seocho-gu (population movement data)

The period with the most dramatic price fluctuations was 2019-2023.

2) 2019~2023 data + infrastructure level 1 (Yongsan-gu) + Yongsan-gu (population movement data)

3) 2013~2023 data + infrastructure level 5 (Seoul) + Seoul (population movement data)





- res  
ult  
Sejong City actual  
transaction index  
forecast
- 1) Coefficient of determination: 0.968
  - 2) Train MSE: 6.3516 / Test MSE: 15.6219
  - 3) Average MSE: 24.9379
- After training by entering data from 13.01 to 23.11,  
Predicting the actual transaction index for December 23
  - **Actual transaction index: 122.0**
  - **Model predicted actual transaction index: 130.33**
- \* Many performance improvements are needed**

1. Apartment price prediction service

Hypothesis: Higher infrastructure levels lead to higher real estate prices

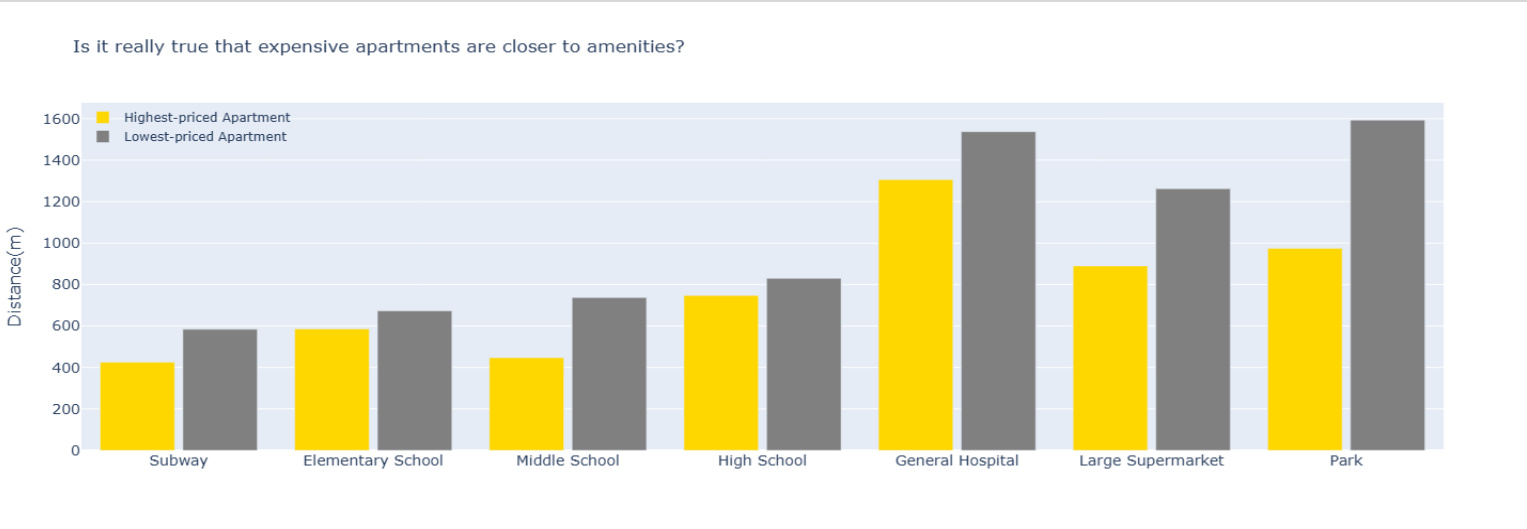
insp

(Scope) Old apartments in Seoul

(Target) Highest price apartment,  
on!  
lowest price apartment

(Confirmation) Distance between  
apartment and infrastructure

Results: The more expensive the apartment, the closer it is to



Calculate the average distance between apartments and



## 2. Future development direction

## Scenario-Based Property Recommendation Solution (Case Study of Virtual Character Kim)



### Standard 1



**Preferred  
Infrastruct**



“I am a working person and I live with my wife.  
I am looking for a place close to the subway station.  
I like to run and run around my house.  
“There is a park and it looks nice.”

### Standard 2

**Homeless, age, asset size, household  
size, income level**

**Non-homeowner, age: 20s-30s, asset  
size: 500 million won**

**Residential capacity: 2 people, Annual  
salary: 50 million won**

development  
direction

Customer-specific data collection and infrastructure  
accessibility classification

infrastructure	1st grade	2nd grade	3rd grade	Grade 4	Grade 5
subway	425.08	464.81	504.54	544.27	584
elementary school	586.36	608.03	629.7	651.37	673.04
middle school	447.76	520.2	592.64	665.08	737.52
high school	747.2	767.99	788.78	809.57	830.36
general hospital	1306.96	1364.48	1422	1479.52	1537.04
hypermarket	890.44	983.55	1076.66	1169.77	1262.88
park	974.44	1129.12	1283.8	1438.48	1593.16

customer	age	Desired amount	Number of people in the household	Preferred Infrastructure - 1st Choice	Preferred Infrastructure - 2nd Choice
Shin Jeong-yoon	20	150,000,000	1	subway	hypermarket
Park Joo-chan	30	300,000,000	3	elementary school	general hospital
Mr. Oh Seung-pil	40	500,000,000	2	hypermarket	subway
Voice actors	50	1,000,000,000	2	high school	park
Mr. Lee Ui-jae	60	650,000,000	4	general hospital	hypermarket
Python	70	450,000,000	2	park	subway
Deep Learning	80	300,000,000	2	general hospital	hypermarket
Pandas	90	2,000,000,000	1	park	hypermarket



End of Page  
- thank you-

