# Portfolio

●●●

Jinyoung Kim

# Contents

- Real-world data handling experience
  - Cleansing and interpolation of electricity usage data (error detection and correction)
  - EDA and statistical approach to energy usage data
  - Machine learning with Kakao Brunch article data (tf-idf, word2vec, doc2vec)
  - Analysis on transport card data and relevant bus routes (QGIS)
- Development of customized programs for productivity
  - Automatic AutoBase capturer
  - Electricity usage data crawler and calculator
  - Electricity data checker (for faster visualization and modification of the data)
  - $CO_2$ emission checker (for $CO_2$ emission trading)
  - Naver news article crawler
- Contest
  - Smart City Simulation Contest by LH(Korea Land and Housing Corporation) - grand prize

# Real-World Data Handling Experience

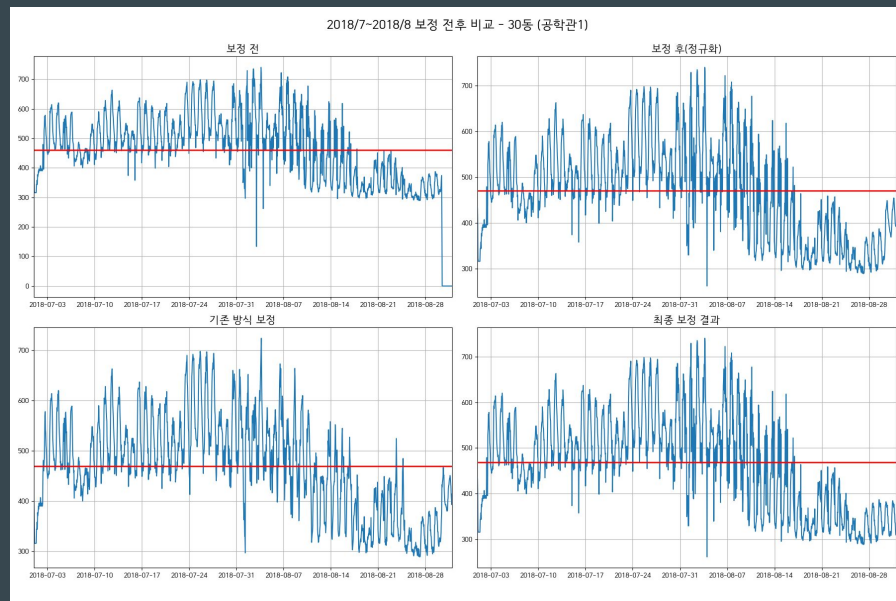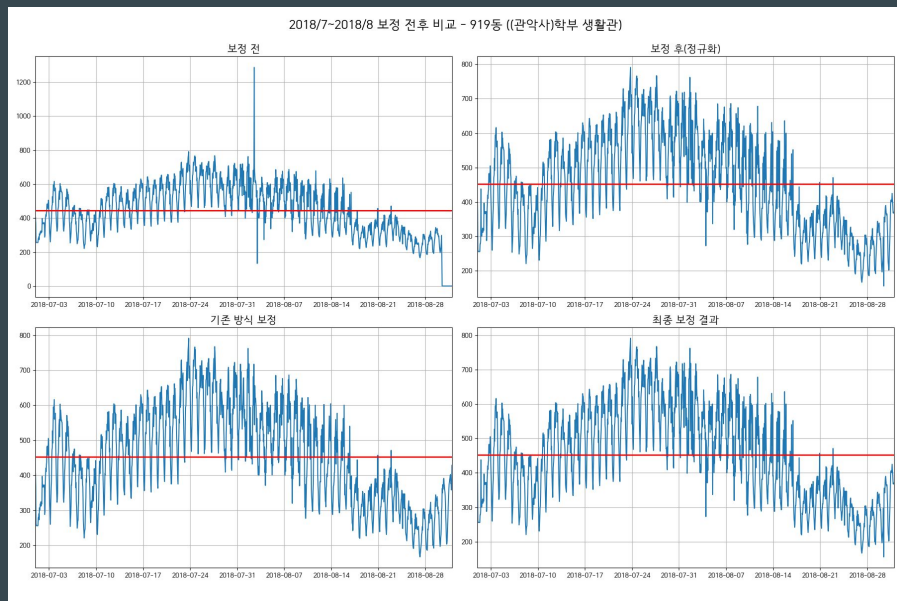# Cleansing/Interpolation of Electricity Usage Data

- Given data
  - Hourly electricity usage data per building (total 156 buildings)
- Types of errors
  - Unrealistically large or negative usage
  - Complete loss of data
  - Repeated usage for too many times
- Reason of errors
  - Malfunction of individual meters
  - Connection loss between meters and the control server
  - Other unknown reasons
- Rate of errors
  - Some 10% with minor fluctuation month by month
  - After the cleansing process, more than 90% of the errors are replaced with plausible values

# Cleansing/Interpolation of Electricity Usage Data

- Error detection strategies
  - (variance of difference) Using the fact that electricity usage tends to change continuously, consider a usage value with a large gap between that of the previous hour as an error.
  - (confidence interval) For every value, sample the values ±1 hour, ±1 day, and ±1~4 years to make a sample group and get the 98% confidence interval from the group. If the value is out of the range, consider it an error.
- Interpolation (correction) strategies
  - (weighted mean) For every erroneous value, sample the values ±1~4 weeks and ±1~2 years to make a sample group, and get a weighted mean as a replacement. weights are given inversely proportional to time distance from the erroneous value.
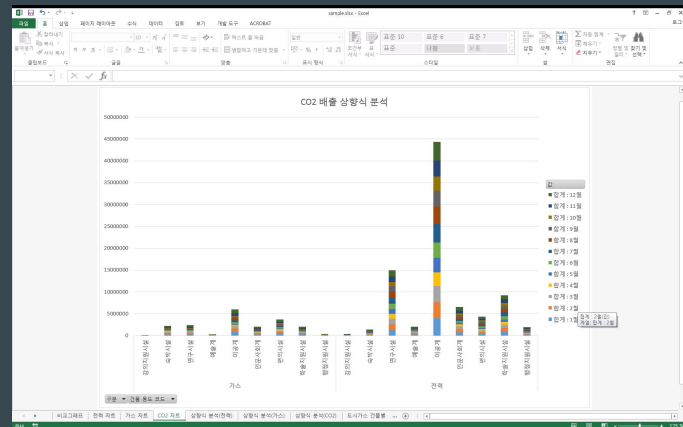  - (mean) Use the mean value or the median of a sample group as a replacement.
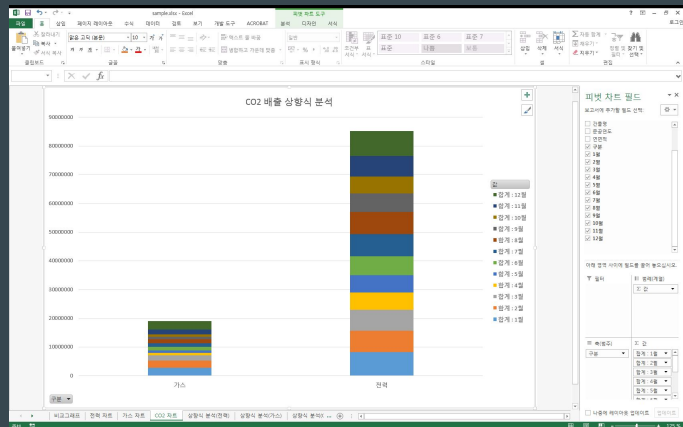
# Cleansing/Interpolation of Electricity Usage Data

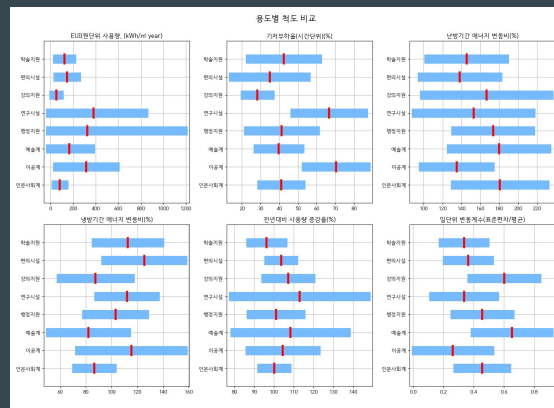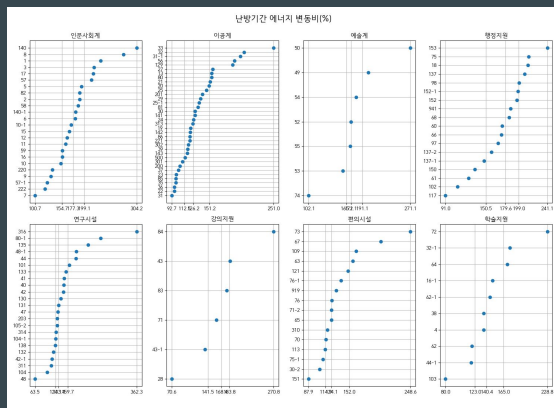● Sample comparison graphs (drawn by python matplotlib)
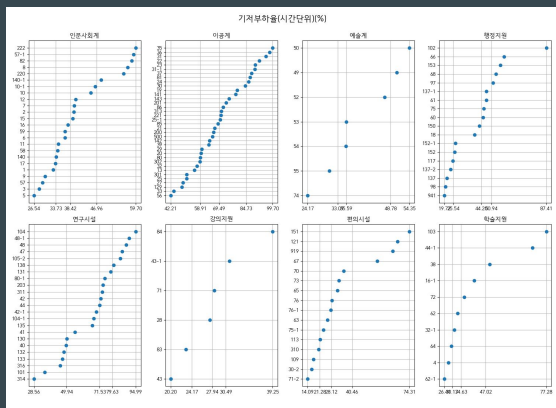
# EDA/Statistical Approach to Energy Usage Data

- CO2 emission status visualization
  - Using hourly electricity usage data per building and hourly natural gas usage data per building, calculate the CO2 emission value.
  - Summarizing the value according to the metadata of the buildings, such as department, purpose, gross floor area, etc.
  - Sample graphs in Excel

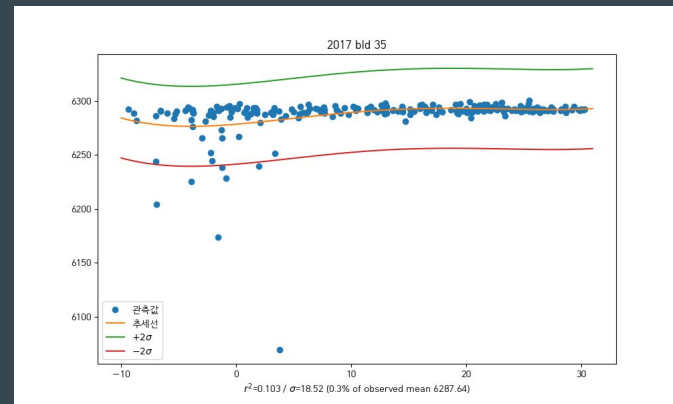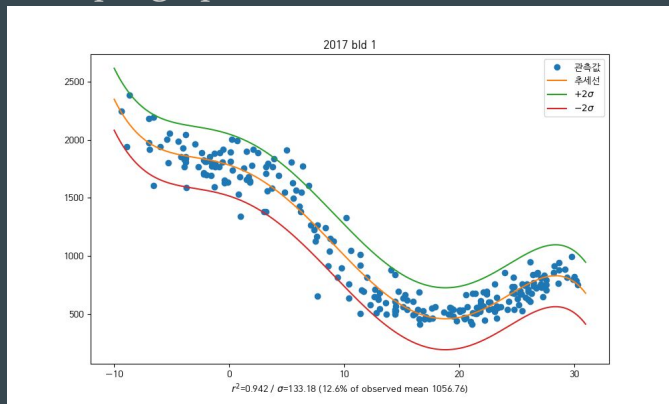# EDA/Statistical Approach to Energy Usage Data

- Characteristic index analysis for each building
    - With the help of metadata of the buildings and utilizing some data-driven methods such as EDA(exploratory data analysis) and PCA(principal component analysis), several indices that represent the characteristics of electricity use of a building were designed.
    - Samples of the result graphs

# EDA/Statistical Approach to Energy Usage Data

- Correlation/regression analysis
  - Goal : to understand the characteristics of patterns of energy usage(especially electricity)
  - Domain Knowledge : electricity usage is always high on summers and winters - temperature can be a major factor deciding the amount of energy use
  - Variables : daily mean temperature and daily electricity usage value per building
  - Regression formula : 5th polynomial (useful to detect inflection points)
  - Sample graphs

# Machine Learning with Kakao Brunch Article Data

- Goals
  - To design a measure to check the similarity among articles to make a recommendation to Brunch users
- Given data
  - A list of Brunch users and their followers
  - A list of Brunch articles and tokenized text of them
  - Access record of articles in time sequence
- Procedures
  - To make a tf-idf matrix using the titles and the snippets of the given articles
  - To use Tensorflow library, implement a doc2Vec machine learning algorithm and obtain a document matrix
  - To utilize cosine similarity when checking the (virtual) distance between documents
  - To make use of the result for recommending relevant articles customized to each user

# Machine Learning with Kakao Brunch Article Data

- Similarity analysis with tf-idf matrices(left)  /  Doc2vec modeling(right)

```python
with open('titleTfIdf.pkl', 'rb') as f:
    tfidfUniMatrix = pickle.load(f)
```

```python
[39]  ▷ M↓
coss2 = tfidfUniMatrix.dot(tfidfUniMatrix[0].transpose()).toarray()
```

```python
[40]  ▷ M↓
simmmm = [(x, y) for x, y in enumerate(coss2)]
```

```python
[65]  ▷ M↓
with open('titleTfIdf.pkl', 'wb') as f :
    pickle.dump(tfidfMatrix, f)
```

```python
[38]  ▷ M↓
simm = []
count = 0
for i in range(tfidfMatrix.shape[0]) :
    coss = tfidfMatrix[0].dot(tfidfMatrix[i].transpose())
    coss = coss.toarray()[0][0]
    simm.append((count, coss))
```

```python
[42]  ▷ M↓
simm2 = []
for count, cos in enumerate(simm):
    simm2.append((count + 1, cos[1]))
```

```python
[41]  ▷ M↓
sortedSimm = sorted(simmmm, key=lambda x : x[1], reverse=True)
```

```python
[81]  ▷ M↓
titles[[x for x, y in newSim[:50]]]
```

```
0                사진으로 옮기기에도 아까운 리치필드 국립공원 세상 어디에도 없는 호주 Top 10
530947            세상에서 가장 큰 모래섬 프레이저 아일랜드 세상 어디에도 없는 호주 Top 10
98265        세상 어디에도 없는 독특한 맛 캥거루고기를 아시나요 세상 어디에도 없는 호주 Top10
435307     살아온 날보다 살아갈 날이 많은 나의 최선의 선택 호주 프롤로그 세상 어디에도 없는...
480514                          착한 남자 세상 어디에도 없는
517521             반듯하게 계획된 도시 캔버라 세상 어디에도 없는 호주 Top10
4271                   화분 적당함이란 세상 어디에도 없는 것일지 몰라
416810              세상 어디에도 없는 절규남자 후지산 여행기 셋
496235                세상 어디에도 없는 값 Korean Gapjil Air
39472        세계 최대의 고래 출몰 허비 베이 고래 관찰 투어 세상 어디에도 없는 호주 Top 10
```

```python
            .format(generation, loss_val, currentepoch, currentindexposition))

# Validation: Print some random words and top 5 related words
if generation % print_valid_every == 0:
    sim = sess.run(similarity, feed_dict=feed_dict)
    for j in range(len(sample_words)):
        valid_word = word_dictionary_rev[indexOf_sample_words[j]]
        top_k = 5 # number of nearest neighbors
        nearest = (-sim[j, :]).argsort()[1:top_k+1]
        log_str = "Nearest to {}:".format(valid_word)
        for k in range(top_k):
            close_word = word_dictionary_rev[nearest[k]]
            log_str = '{} {},'.format(log_str, close_word)
        print(log_str)

# Save dictionary + embeddings
if generation % save_embeddings_every == 0:
    # Save vocabulary dictionary
    with open(os.path.join(data_folder_name,'data0analyzed.pkl'), 'wb') as f:
        pickle.dump(word_dictionary, f)

    # Save embeddings
    model_checkpoint_path = os.path.join(os.getcwd(), data_folder_name, 'word2vec_embeddings.ckpt')
    save_path = saver.save(sess, model_checkpoint_path)
    print('Model saved in file: {}'.format(save_path))
```
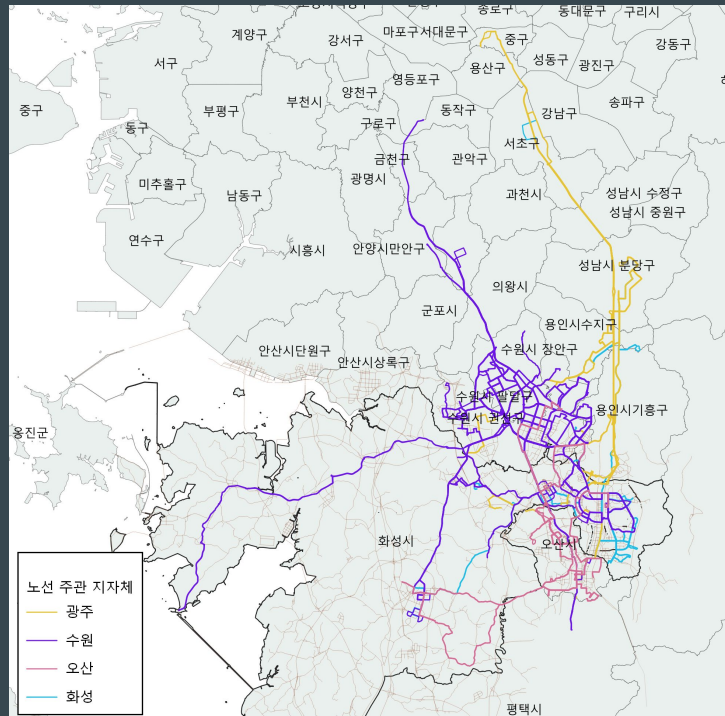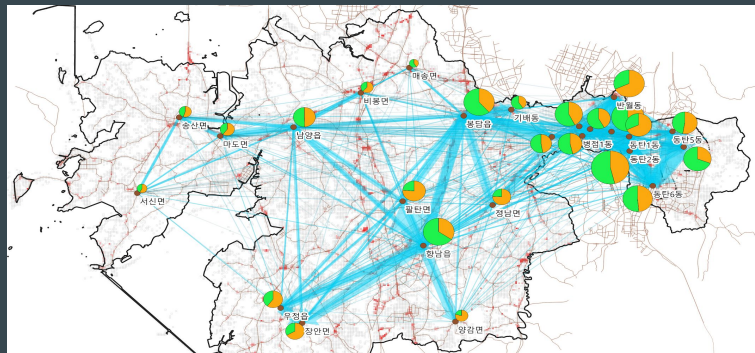
```
Creating Dictionary
Creating Model
Starting Training
Loss at step 1000 : 272.0592956542969 / current epochs : 0 / current sentence : 108
Loss at step 2000 : 242.02508544921875 / current epochs : 0 / current sentence : 230
Loss at step 3000 : 221.26219177246094 / current epochs : 0 / current sentence : 334
Loss at step 4000 : 201.0927276611328 / current epochs : 0 / current sentence : 458
Loss at step 5000 : 182.1846160888672 / current epochs : 0 / current sentence : 585
Nearest to 19517/NNG: 2828/NNG, 19389/NNG, 3292/NNP, 4127/NNG, 11190/NNG,
Nearest to 9822/NNG: 8792/NNG, 2653/XR, 367/NNG, 17798/NNG, 4883/NNG,
Nearest to 1400/NNG: 4847/NNG, 7349/NNG, 1705/VV, 2740/NNG, 1408/NNG,
Nearest to 16277/NNG: 11830/NNG, 30581/NNG, 5253/NNG, 19067/NNG, 10961/VV,
Model saved in file: E:\kakao\temp3\word2vec_embeddings.ckpt
Loss at step 6000 : 183.94854736328125 / current epochs : 0 / current sentence : 711
Loss at step 7000 : 175.04400634765625 / current epochs : 0 / current sentence : 841
Loss at step 8000 : 163.82937622070312 / current epochs : 0 / current sentence : 958
```

# Analysis on Transport Card Data and Relevant Bus Routes

- Goals
    - To examine the status of bus service in Hwaseong, a city in Gyeonggi Province
    - To present realistic methods to improve the current system
- Given data
    - Geojson file of the road network in Hwaseong and its vicinity
    - Geojson file of the distribution of floating population and residents in Hwaseong
    - Transport card use data of 1st ~ 4th of July, 2018
    - Information on the current bus routes operating in Hwaseong and Suwon (incomplete)
- Procedures
    - Parse the given geojson files to make them easier to manipulate in QGIS
    - Obtain position data of the bus routes from Naver Map
    - Pre-process the transport card use data (removing rows with missing values, dropping columns unnecessary for analysis, etc)
    - Overlap different kinds of the shape files in QGIS to get an insight

# Analysis on Transport Card Data and Relevant Bus Routes

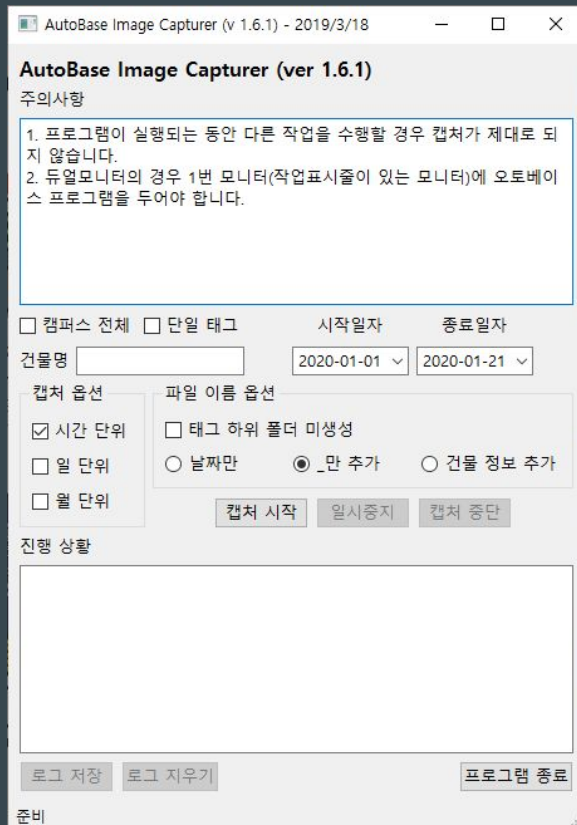- Population density and flow analysis(left) / Bus route analysis(right)

# Development of Customized Programs

# Automatic AutoBase Capturer

- **What is AutoBase?**
  - An operating program of electricity meters (in this case, meters installed in the school campus)
  - An end user can see information on each meter and check the usage, but cannot download it
- **Goals**
  - To save screens of the program showing electricity usage without human intervention
- **Conditions**
  - AutoBase only works under the intranet of the campus
  - Buildings with meters are shown on the campus map, and an end user can click the building to see the meters. The meters should be clicked to see the usage data.
  - Other members of the office do not have any knowledge on programming
- **Strategies**
  - To use PyQt5 library to make a GUI program with Python
  - For ease of development, the program gets the whole control of a mouse and a keyboard of a computer when it is working

# Automatic AutoBase Capturer

- Program interface
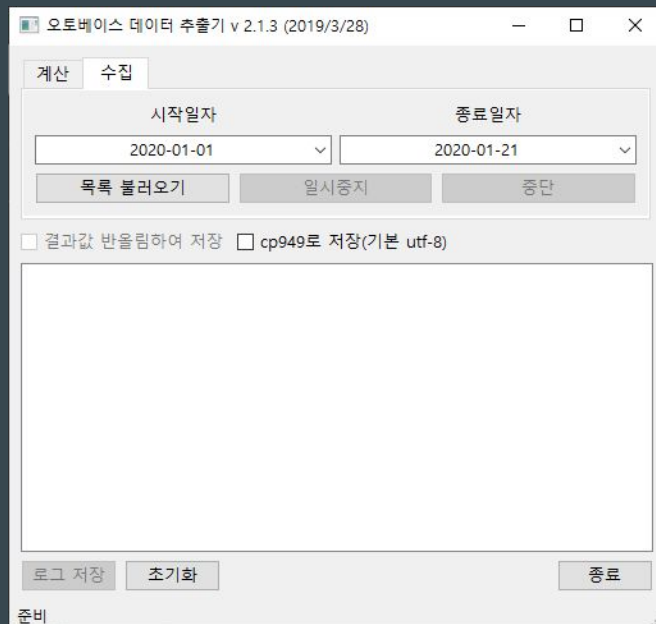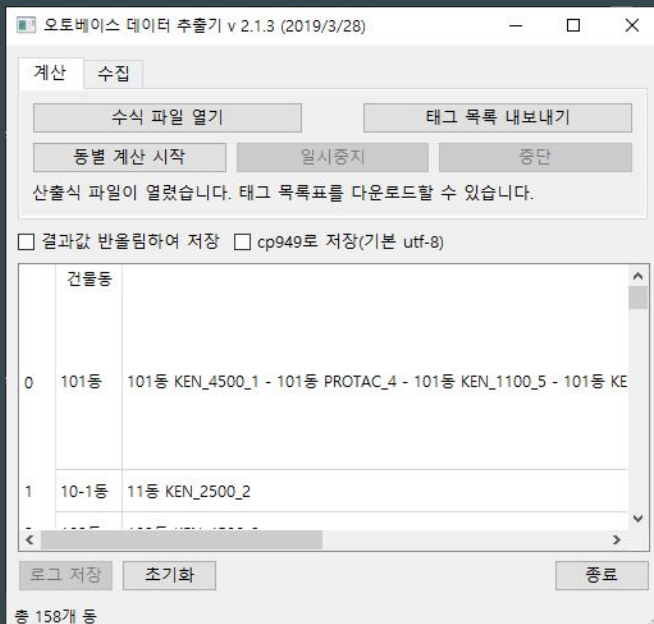  - Video of actual use → https://youtu.be/zM8Vzt6fT-8

# Electricity Usage Data Crawler and Calculator

- Issues
  - AutoBase does show the usage values, but it does not have any function to download the data.
  - It is often needed to calculate the values from different meters to get the usage of a single building
- Goals
  - Automatically crawl the electricity usage data on-demand
  - Calculate the meter-wise data to obtain building-wise usage data according to the given formulae
- Conditions
  - AutoBase gets the usage data through http from the server.
- Strategies
  - With a web debugger (Fiddler in this case), find the packet that AutoBase sends to request the required data and mimic the pattern with Python Requests library
  - Use a single format for saving the usage values on meters and buildings alike
  - To use PyQt5 library to make a GUI program with Python

# Electricity Usage Data Crawler and Calculator

- Program interface
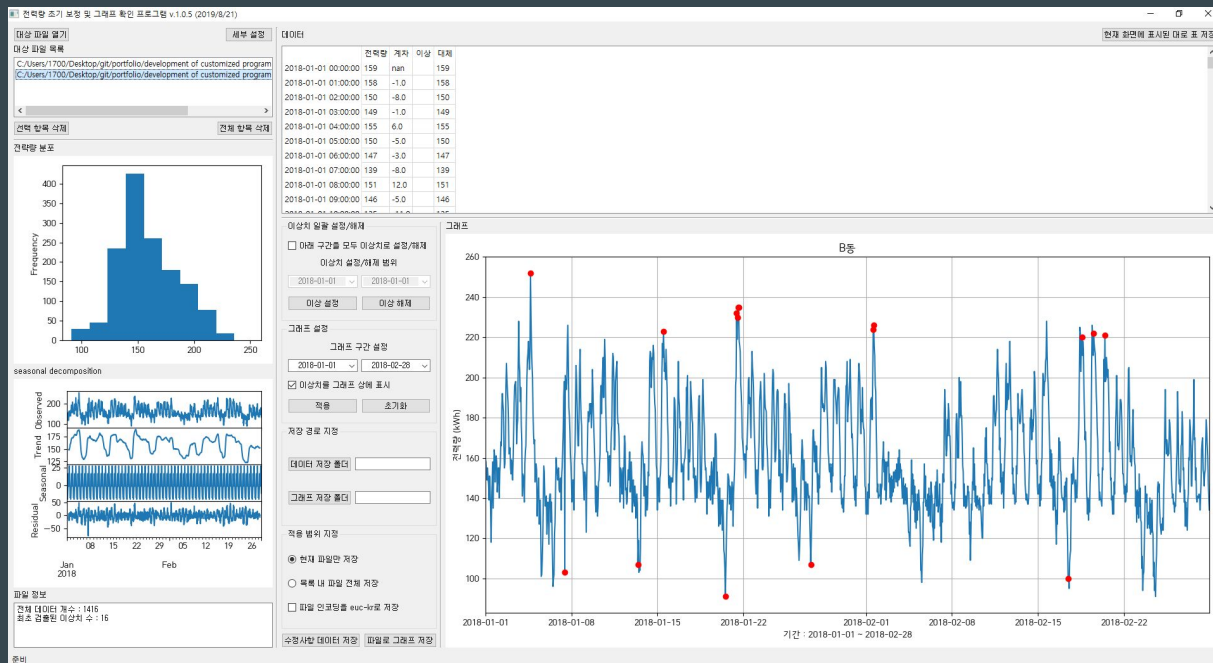  - Video of actual use → https://youtu.be/AjSUMGZaC8A

# Electricity Data Checker

- Issues
  - Plotting graphs of hourly electricity usage data gets trickier as the data cumulates.
  - Though data cleansing process is done automatically, some of errors are still undetected.
- Goals
  - To provide a fast and convenient tool for plotting graphs of the usage data and replace undetected but invalid values with a zero, which is always detected as an error during cleansing process.
- Strategies
  - To use PyQt5 library to make a GUI program with Python
  - To implement seasonal decomposition, a method used for time series analysis, to aid the examination of possible errors
  - To allow a user to choose which encoding to use when saving the result as a csv file, considering that Excel is the main program in the center. (csv and text files saved with Unicode are tricky to open with Excel)

# Electricity Data Checker

- Program interface
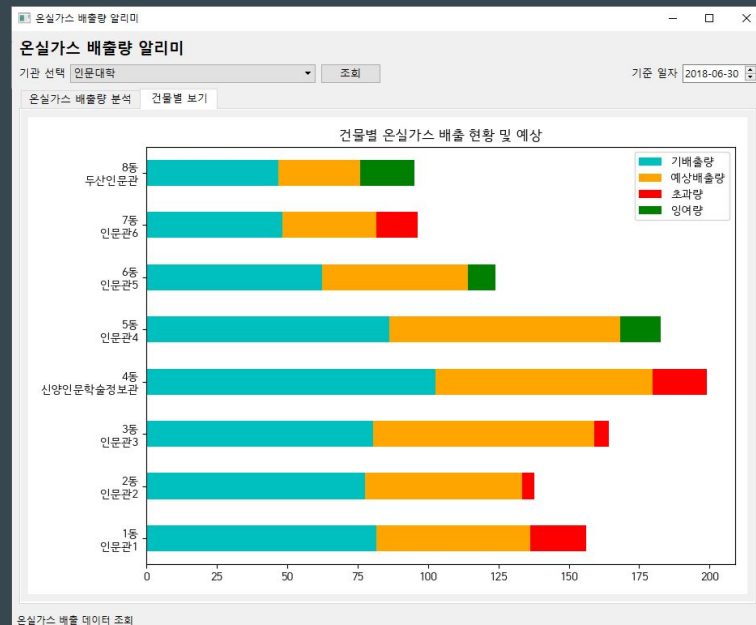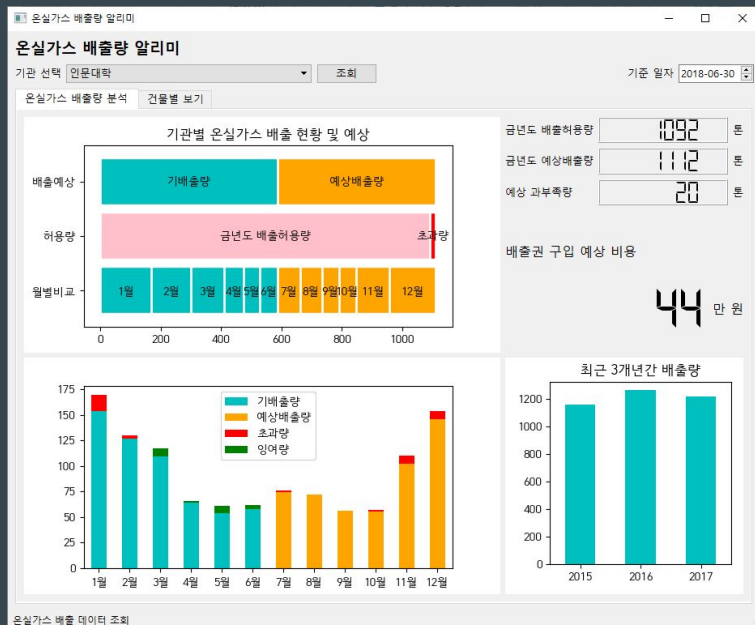  - Video of actual use → https://youtu.be/s6ttk3NJbQs

# CO2 Emission Checker

- Issues
  - Due to the advent of emission trading, the necessity of notifying the amount of CO2 emissions to the administrators of the departments in the campus.
- Goals
  - To make an easy-to-understand program showing the CO2 emission per department
- Conditions
  - CO2 emission by electricity and natural gas must be combined properly
  - CO2 emission should be shown per department and per building when necessary
  - Distributable version of the program is to be developed by an outside company
- Strategies
  - To make a prototype of the program with Python to review its composition
  - To provide the prototype to the outside software company for reference

# CO2 Emission Checker

- Program interface
  - Video of actual use → https://youtu.be/F09TqxCIGA4

# Naver News Article Crawler

- Issues
  - While participating in a research group, the need to construct database on conflicts among residents on property development was raised.
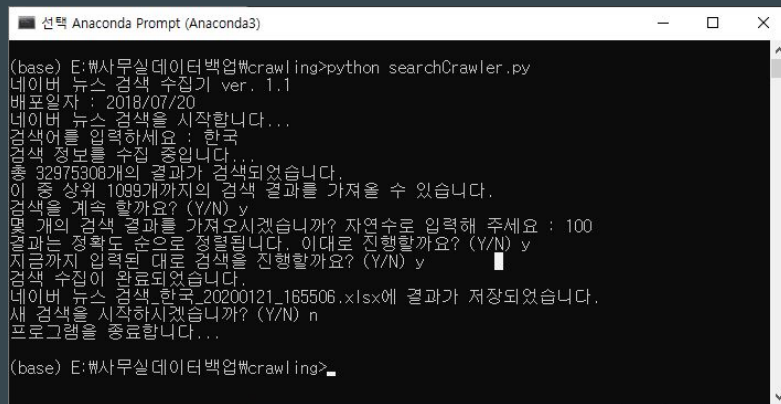- Goals
  - To gather news articles of a topic on Naver News
- Strategies
  - To utilize an open API of news search provided for developers by Naver
  - To make a CUI program to shorten the development period
  - To package the Python code file into an exe file with Pyinstaller library
- Paper
  - A study on conflict DB construction using python-based web-crawling (2018, participated as a co-author)

# Contest

# Smart City Simulation Contest

- Overview of the contest
  - Held in 2018 by LH(Korea Land and Housing Corporation)
  - Participants were expected to freely design 5-1 district of Sejong City under the basic concept and development plan of the place and represent it with Cities: Skylines, a city-builder simulation game.
  - The grand prize winner was expected to revise their work according to requests from the host and submit a 5-minute complete video clip of their work on Cities: Skylines
- Participation progress and result
  - Participated as a team of two (major role of mine was to construct the city on the game and edit the video clip)
  - Title of work was "Baraon(바라온)" and its major concept was "An Innovative Space Planning Resembling Artificial Neural Network Structure"
  - Awarded the grand prize with 10,000,000 KRW as a reward
  - Conducted the following work until late 2019