

에너지 빅데이터 활용 데이터 사이언스 경진대회

유 광 남

데이터 파일

1. Train.csv

> 1300 세대의 2016.07.26 ~ 2018.06.30 전력 데이터

2. Test.csv

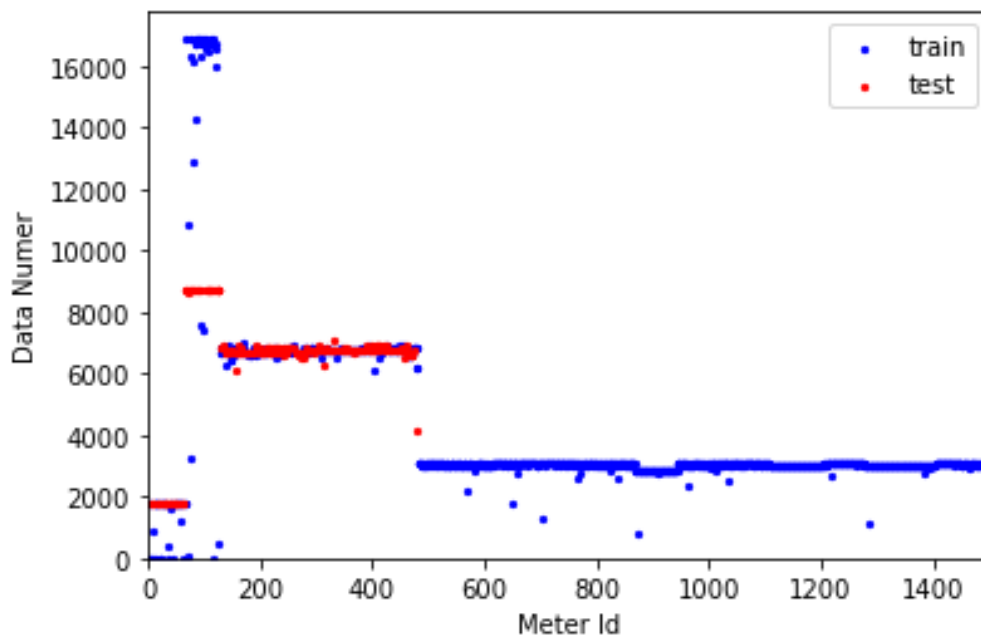
> 200 세대의 2017.07.01 ~ 2018.06.30 전력 데이터

3. 인천_시간별_기상자료(16-18)_축소_7월1일.csv (주최측 제공)

> 인천 지역의 2016.07.20 ~ 2018.07.01 날씨 데이터

결측치

1. Train.csv, Test.csv 파일 결측치로 데이터 개수에 차이가 있음
2. ID 481번 기준으로 데이터 개수의 분포가 다름
> 0 ~ 481 번 데이터만 사용하는 것을 고려

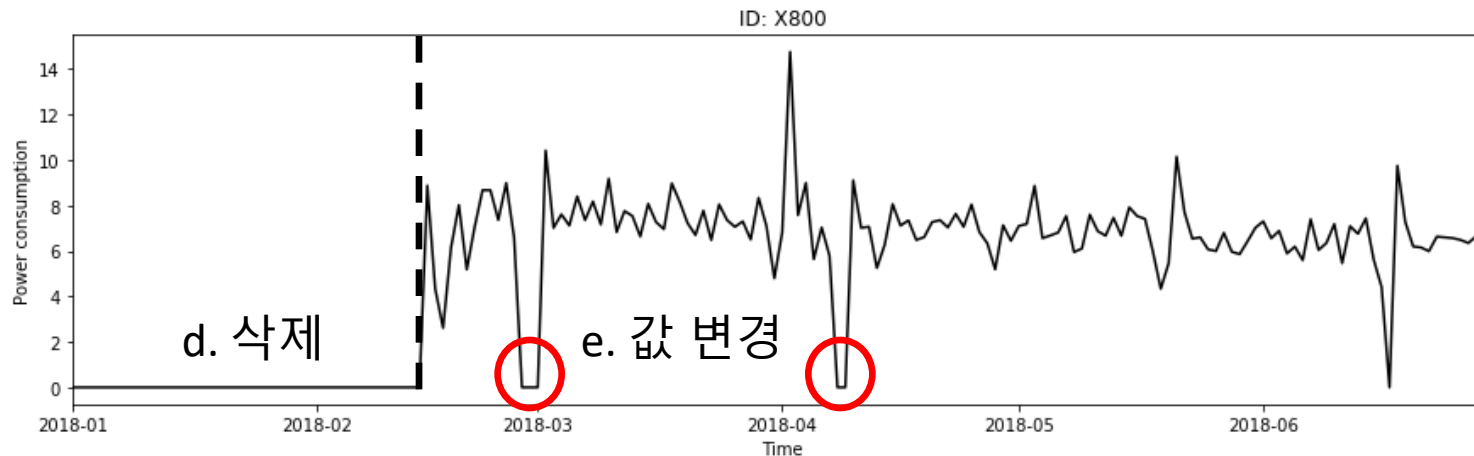


결측치 처리

대회 참가하기



- a. 결측치는 모두 0으로 변경
- b. 0.003 보다 작은 값은 0으로 변경
- c. 시간별 전력을 하루 단위로 합쳐 날짜별 전력으로 변경
- d. 유효한 데이터만 사용하기 위해 앞뒤로 0인 구간을 삭제
- e. 데이터 중간에 0이 있는 경우 값 변경 (ex. 평균값)



Feature

대회 참가하기



1. 전력 데이터 정규화 (*meter id 별로 따로 처리)
> 최대값으로 나눔 > 로그를 씌움 > 평균으로 뺌 > 표준편차로 나눔
2. 전력데이터에 해당하는 날짜의 시간, 요일 정보를 사용
> 시간은 주기가 365일, 최대값이 8월 1일인 사인함수를 사용
3. 2018.07.01 을 예측하는 경우 날씨 데이터 사용을 고려함
4. 2018.07.01 이후를 예측하는 경우 날씨 데이터 사용 x
5. 예측하는 날의 약 10주 전 데이터를 모아 1차원 벡터로 사용

모델

대회 참가하기



LightGBM

조건에 따라 18개 모델을 생성함

> ID: (0~481 / 0~1500)

> 전력데이터 중 0에 대한 처리를 3 종류로 분류

> 2018.07.01 예측에 날씨 데이터 (사용 / 사용 X)

> 기상 모델 사용시 2018.07.02 예측을 위해 (74 / 75)일 데이터 사용

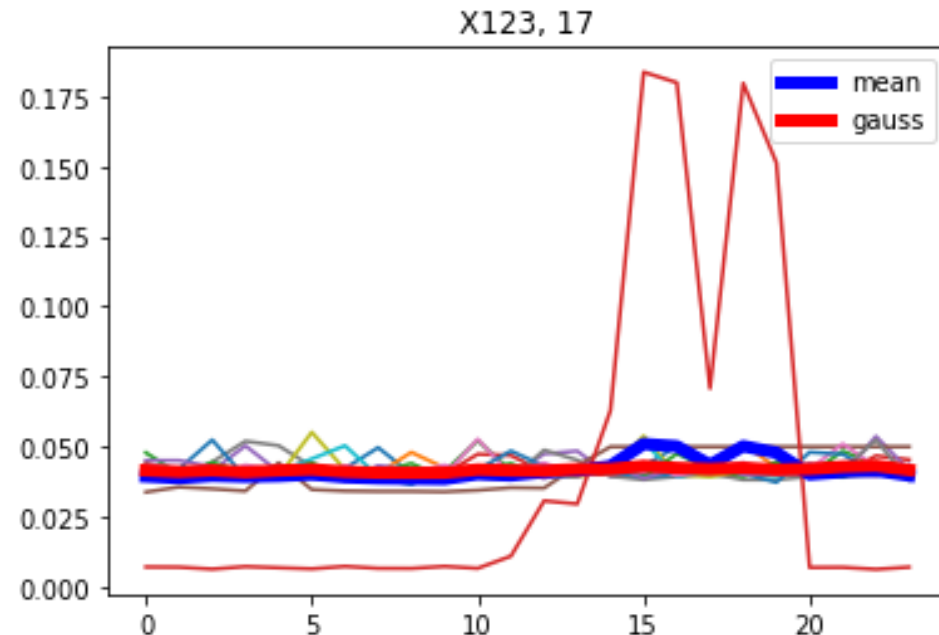
시간별 전력 수요 예측(18.07.01)

대회 참가하기



2018.07.01의 시간별 수요 예측을 위해 2018.03.01~2018.06.30 기간
내 일요일에 해당하는 시간별 수요의 분포를 구함

Gauss 함수로 weight를 주어 좀 더 일반적인 대표값을 구함

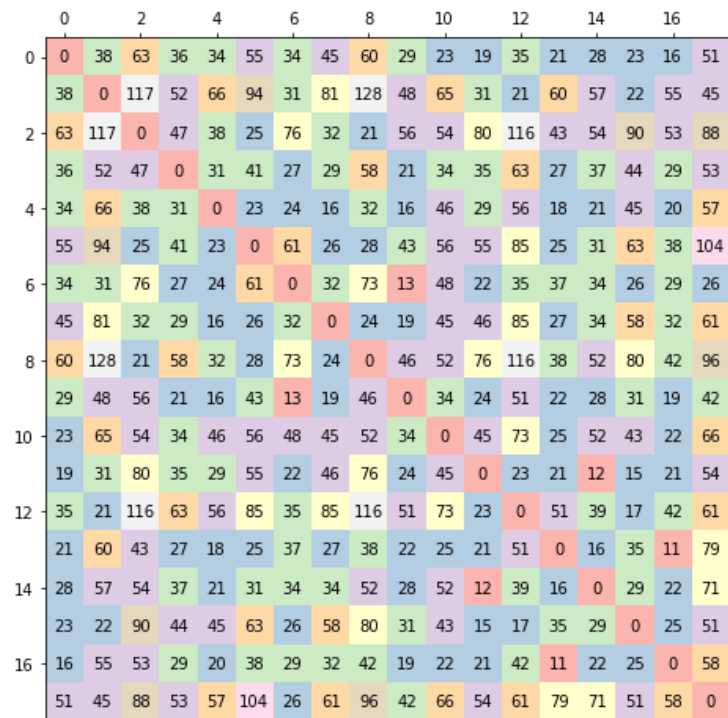


모델 선택

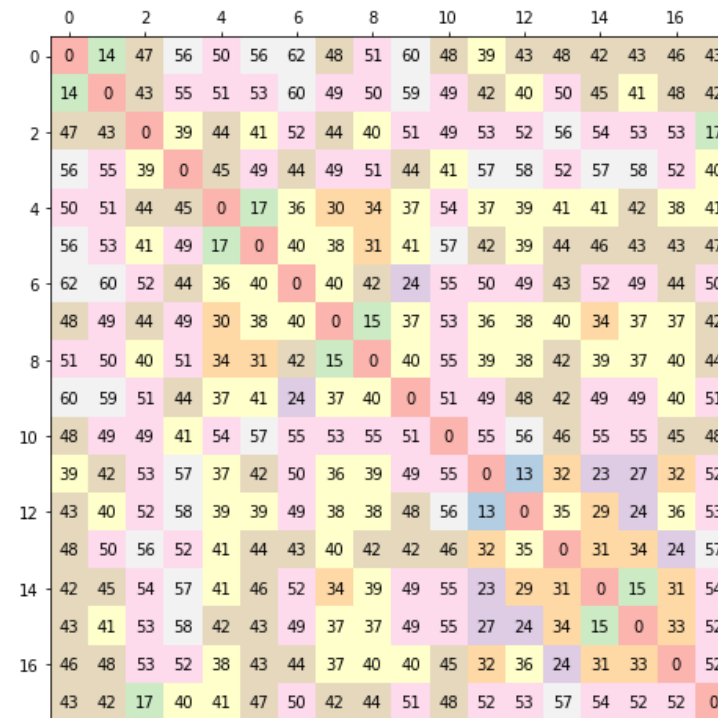
대회 참가하기



1-R2 score



SMAPE

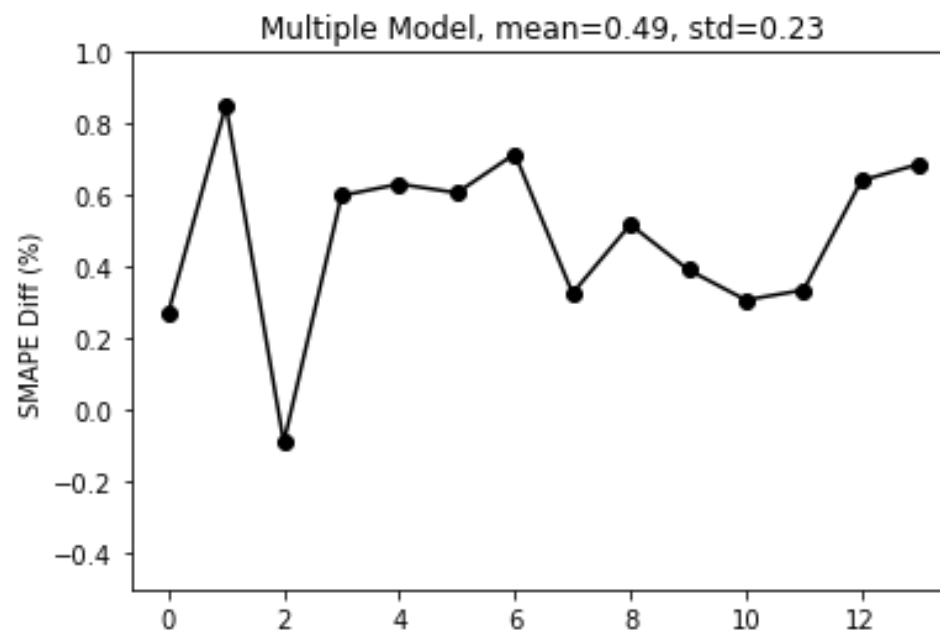
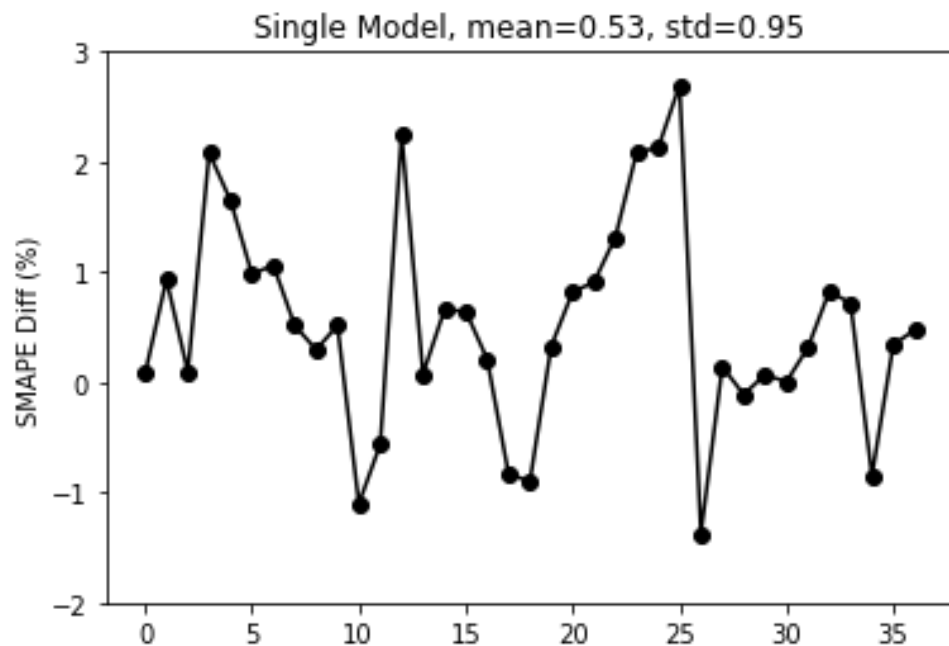


18개 결과에 대해 1-R2 score, SMAPE를 계산, 각 score의 중간값 보다 작은 데이터만 선별하고 평균을 구함

<https://dacon.io>

모델 선택

대회 참가하기

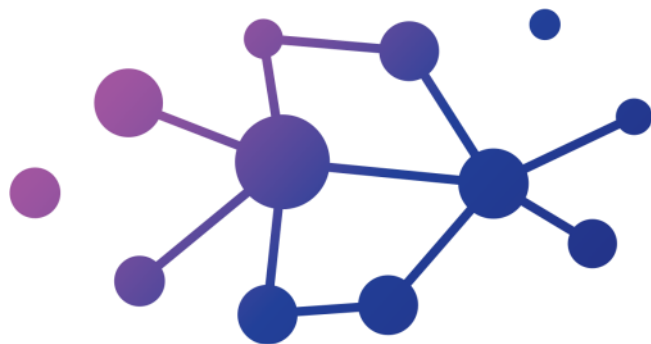


Single model, Multiple model 의 private, public score 차이 (백분율)

Single model: mean=0.53, std=0.95, private score가 public score에 비해 약 0.5% 높다.

Multiple model: mean=0.49, std=0.23, single model 처럼 private score가 public score에 비해 약 0.5% 높음, single model에 비해 std가 4배 작아 더 안정적으로 예측가능함

THANK YOU



대회 참가해보기