

2023.09.17

LG aimers 10등 solution

TEAM 치킨마요덮밥

목차

01 EDA 및 데이터 전처리

02 모델 학습

03 추론

01. EDA 및 데이터 전처리

01 EDA 및 데이터 전처리

1) train 데이터 특징

28894개의 고유ID의 480일치 판매량으로 이루어진 데이터.

	ID	제품	대분류	중분류	소분류	브랜드	쇼핑몰	2022-01-01	2022-01-02	...	2023-04-23	2023-04-24
0	0	B002-00001-00001	B002-C001-0002	B002-C002-0007	B002-C003-0038	B002-00001	S001-00001	0	0		0	0
1	1	B002-00002-00001	B002-C001-0003	B002-C002-0008	B002-C003-0044	B002-00002	S001-00001	0	0		0	0
2	2	B002-00002-00002	B002-C001-0003	B002-C002-0008	B002-C003-0044	B002-00002	S001-00001	0	0		0	0
3	3	B002-00002-00003	B002-C001-0003	B002-C002-0008	B002-C003-0042	B002-00002	S001-00001	0	0		0	0
4	4	B002-00003-00001	B002-C003-0042	B002-C002-0001	B002-C003-0003	B002-00003	S001-00010	0	0		0	8

⋮

01 EDA 및 데이터 전처리

2) 데이터 정제(LSTM)

판매량이 일정 기간 연속으로 0인 구간을 확인 후
이상 요인(품질, 공급부재, 미등록 등)의 이유로 판단되어 제거

	ID	제품		2022-01-01		2023-04-23	2023-04-24
0	0	B002-00001-00001	...	0	...	0	0
1	1	B002-00002-00001		0		0	0
2	2	B002-00002-00002		0		0	0
3	3	B002-00002-00003		0		0	0
4	4	B002-00003-00001		0		0	8

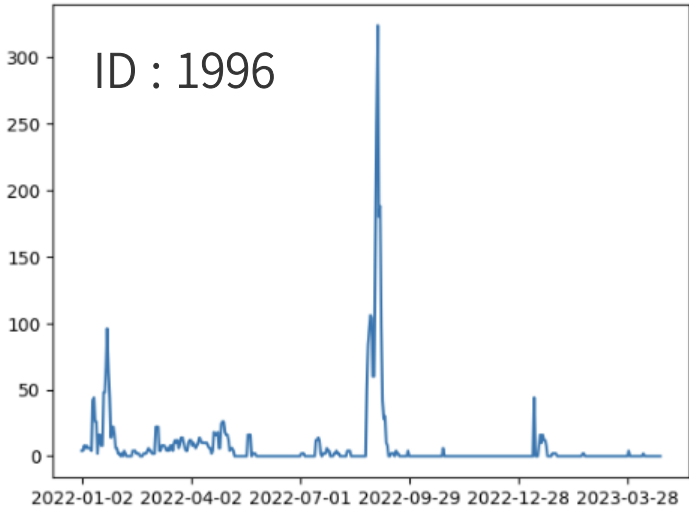
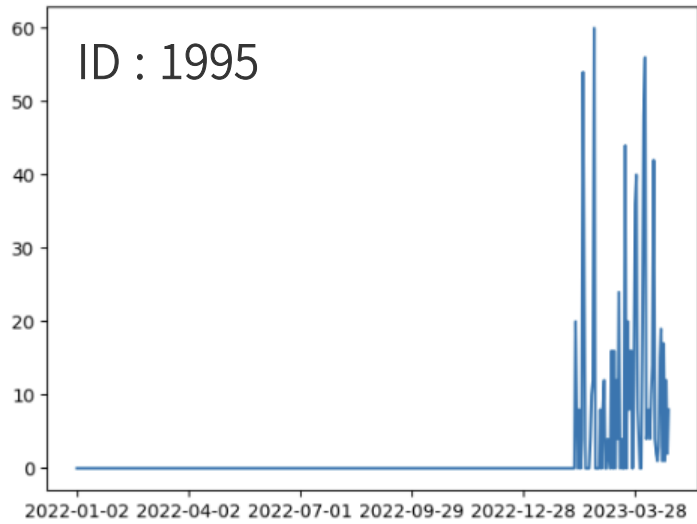
⋮

판매량 30일 이상 0인 구간 제거

```
## 판매량이 30일 이상 연속으로 0인 구간 확인
fin_list = []
for k in tqdm(range(train_data.shape[0])): # 28894개에 대해 반복
    a = train_data.iloc[k, 4:] # 각 제품에 대해
    period = CFG['PERIOD'] # 0인 기간 30
    fir_list = [] # 제품별 삭제 구간 리스트
    i = 0
    while i < len(a)-1: # 전체 길이 동안
        sec_list = [] # 구간 리스트
        if a[i] == 0: # 판매량 값이 0이면
            sec_list.append(i) # 판매량 0 시작 지점
            while a[i] == 0: # 판매량이 0인 동안
                i += 1
            if i == 479:
                break
            sec_list.append(i-1) # 판매량 0 종료 지점
            if (sec_list[1] - sec_list[0]) >= period: # 판매량이 일정 구간 이상
                #이면(30일)
                fir_list.append(sec_list) # 삭제 구간 리스트에 추가
            i+=1
    fin_list.append(fir_list) # 각 제품별 삭제구간 리스트에 i번째 제품 삭제
    # 구간 추가
```

01 EDA 및 데이터 전처리

2) 데이터 정제(ML)



ID별 판매량 차이가 커 ID별로 예측하는 것이 예측력을 높일 수 있다 판단
데이터셋 형태를 다음과 같이 변환

	ID	COUNT	제품	대분류	...	Year	Month	Day	...	COUNTS_56	COUNTS_57	COUNTS_58	COUNTS_59
0	0	0	B002-00001-00001	B002-C001-0002		2022	3	1		0	0	0	0
1	0	0	B002-00001-00001	B002-C001-0002		2022	3	2		0	0	0	0
2	0	0	B002-00001-00001	B002-C001-0002		2022	3	3		0	0	0	0
3	0	0	B002-00001-00001	B002-C001-0002		2022	3	4		0	0	0	0
4	0	0	B002-00001-00001	B002-C001-0002		2022	3	5		0	0	0	0
⋮													
5783959	15889	0	B002-03799-00010	B002-C001-0002		2023	4	4		0	0	0	0

01 EDA 및 데이터 전처리

3) 파생변수 생성

(1) 주말, 휴일, 공휴일 변수 생성 (ALL)

- 주말 및 공휴일이 판매량에 영향을 미칠 것으로 예상하여 주말, 공휴일, 휴일 변수 생성

	Year	Month	Day	Weekday	Weekday_Name	주말여부	공휴일여부	휴일여부
2022-01-01	2022	1	1	5	토	Y	Y	Y
2022-01-02	2022	1	2	6	일	Y	N	Y
2022-01-03	2022	1	3	0	월	N	N	N
2022-01-04	2022	1	4	1	화	N	N	N
2022-01-05	2022	1	5	2	수	N	N	N

01 EDA 및 데이터 전처리

3) 파생변수 생성

(2) 판매량 변동량 역수 변수 생성(LSTM)

- 판매량의 급한 변동에 대한 학습효과를 완화시키기 위해 전일 대비 판매량 변동량의 역수를 취한 변수를 생성

$$V2_t = \begin{cases} \frac{1}{y_t - y_{t-1}}, & t : 2022 - 01 - 02 \text{ 이후} \\ 0, & t : 2022 - 01 - 01 \end{cases}$$

02. 모델 학습

02 모델 학습

모델 검증 방법

- Rolling window 방법 이용(ML : size = 60, LSTM : size = 75)

모델 설정

- 모델 :

LSTM(Batch_size = 8192, 변동율 변수 추가),

LSTM(Batch_size = 2048, 변동율 변수 추가),

LSTM(Batch_size = 2048, 변동율 변수 삭제)

XGBOOST(n_estimator=200),

RandomForest(n_estimator=200) 사용

02 모델 학습

최종 모델 Parameter

LSTM

LSTM1

- TRAIN_WINDOW_SIZE : 75
- EPOCHS : 10
- LEARNING_RATE : 0.0001
- BATCH_SIZE : 8192
- SEED : 9909
- PERIOD : 30

LSTM2,3

- TRAIN_WINDOW_SIZE : 75
- EPOCHS : 10
- LEARNING_RATE : 0.0001
- BATCH_SIZE : 2048
- SEED : 9909
- PERIOD : 30

ML

- train_window_size : 60

RandomForest

- n_estimators = 200, random_state = 41

XGBRegressor

- n_estimators = 200, random_state = 41

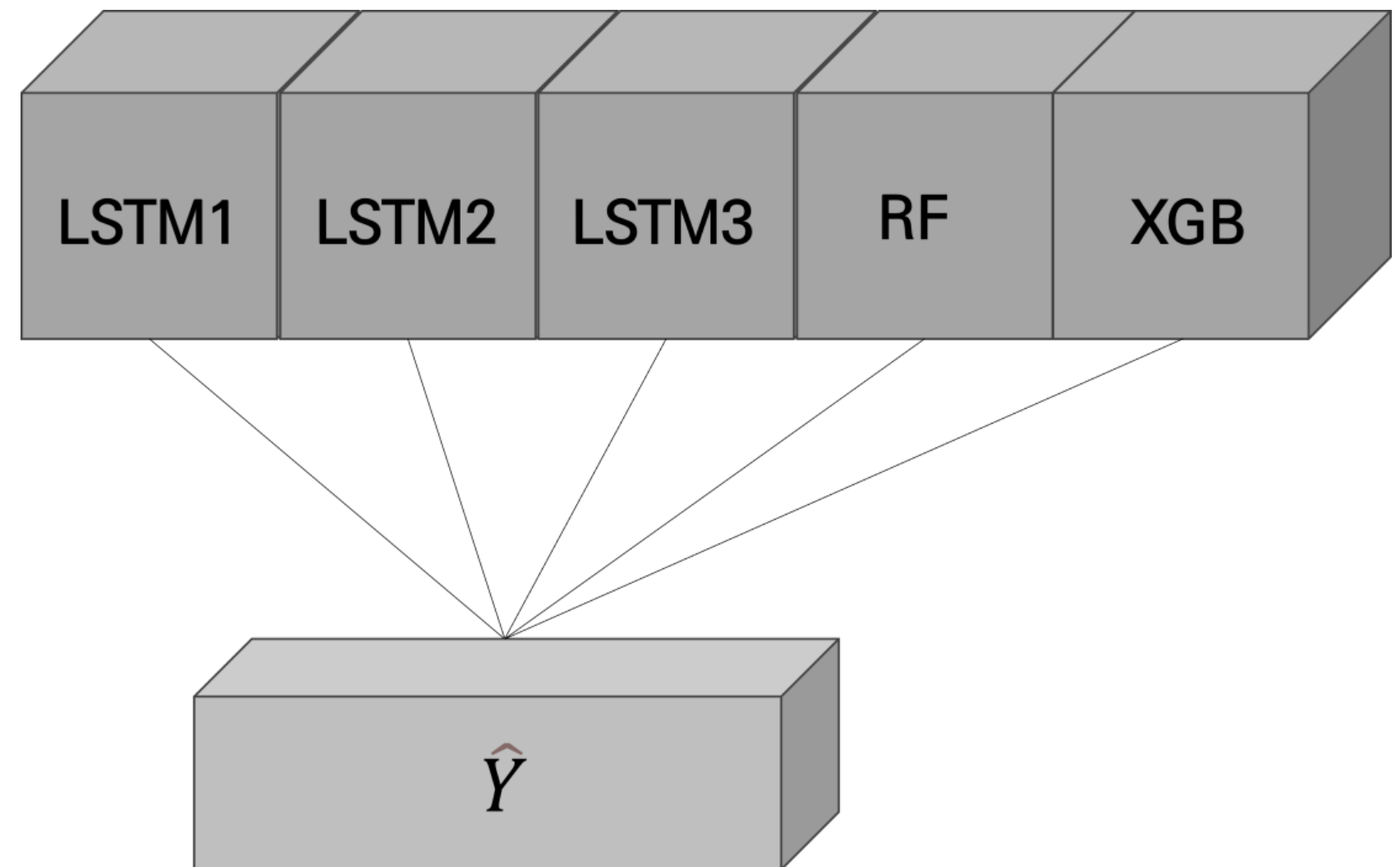
03. 추론

03 추론

예측

5개의 모델(LSTM(1,2,3), XGB, RF)을 앙상블
이후 5개 모델의 예측값에 가중치를 주어 계산함
반올림하여 정수 형태로 반환

LSTM1	LSTM2	LSTM3	XGB	RF
1/6	1/6	1/6	1/4	1/4



03 추론

평가 지표

최소한의 예측 가능성을 확보하고자 최소 판매량을 1로 설정

대분류 별 Pseudo 예측 정확도: $PSFA_m = 1 - \frac{1}{n} \sum_{day=1}^n \sum_{i=1}^N \left(\underbrace{\left(\frac{|y_i^{day} - p_i^{day}|}{\max(y_i^{day}, p_i^{day})} \right)}_{\text{오차}} \times \underbrace{\frac{y_i^{day}}{\sum_{i=1}^N y_i^{day}}}_{\text{(판매)비중}} \right)$

- m: 대분류 index
- i: (대분류 내에서) 제품 index
- y_i^{day} : i번째 제품의 day일의 판매량
- p_i^{day} : i번째 제품의 day일의 예측량

전체 Pseudo 예측 정확도: $PSFA = \frac{1}{M} \sum_{m=1}^M PSFA_m$

Thank you.
