#BIG_DATA_PROGRAMMING

Solar - Prediction

#강예진 60201658 #박선진 60180520 #장유진 60201703





목차

A table of Contents

#1. 데이터 설명

#2. EDA

#3. 모델링

#4. 결과 분석



데이터 설명

Solar Radiation Prediction

- NASA HI-SEAS Project
- 화성과 비슷한 환경의 Hawaii 마우나로아 화산
- 2016년 9월 ~ 12월 동안 측정한 data
- 30000개 정도 데이터
- Kaggle: Solar Radiation Prediction



Part 1 데이터 설명

SolarPrediction.csv (2.96 MB) 坐 # < Detail Compact Column 11 of 11 columns V ☐ Time # UNIXTime Data # Radiation # Temperature # Pressure # Humidity # WindDirection(De... = # Speed TimeSunRise ☐ TimeSunSet 1.6k 34 71 30.2 30.6 103 0.09 360 0 40.5 21Sep19 1.47b 1.48b 1Sep16 31Dec16 20Sep19 21Sep19 1.11 20Sep19 20Sep19 1475228421 9/29/2016 12:00:00 23:40:21 1.21 48 30.46 60 137.71 3.37 06:13:00 18:13:00 1475228124 9/29/2016 12:00:00 23:35:24 1.17 48 30.46 62 104.95 5.62 06:13:00 18:13:00 1475227824 9/29/2016 12:00:00 23:30:24 1.21 48 30.46 64 120.2 5.62 06:13:00 18:13:00 1475227519 9/29/2016 12:00:00 23:25:19 1.2 49 30.46 72 112.45 6.75 06:13:00 18:13:00 1475227222 9/29/2016 12:00:00 23:20:22 1.24 49 30.46 71 122.97 5.62 06:13:00 18:13:00

- UNIXTime >> 유닉스 타임

9/29/2016 12:00:00

23:15:22

1.23

- Date >> 일자

1475226922

- Time >> 하와이 기준 시간
- Radiation >> 태양 복사

- Temperature >>기온

30.46

80

- Pressure >> 기압

49

- Humidity >> 습도
- WindDirection >> 풍향

- Speed >> 풍속

101.18

4.5

- TimeSunRise >> 일출 시간

06:13:00

- TimeSunSet >> 일몰 시간

18:13:00

Part 1 데이터 설명

# UNIXTime =	🗖 Data 🖃	□ Time =	TimeSunRi =	TimeSunSet =
1475229326	9/29/2016 12:00:00 AM	23:55:26	06:13:00	18:13:00
1475229023	9/29/2016 12:00:00 AM	23:50:23	06:13:00	18:13:00
1475228726	9/29/2016 12:00:00 AM	23:45:26	06:13:00	18:13:00
1475228421	9/29/2016 12:00:00 AM	23:40:21	06:13:00	18:13:00
1475228124	9/29/2016 12:00:00 AM	23:35:24	06:13:00	18:13:00



Time
23:55:26
23:50:23
23:45:26
23:40:21
23:35:24

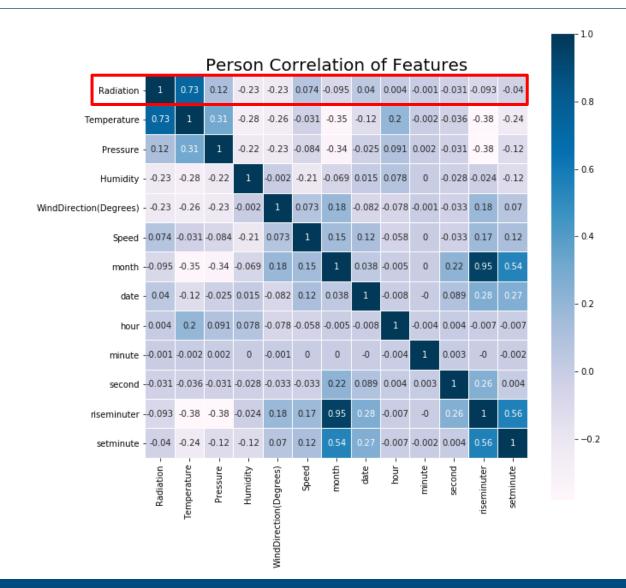
hour	minute	second
23	55	26
23	50	23
23	45	26
23	40	21
23	35	24

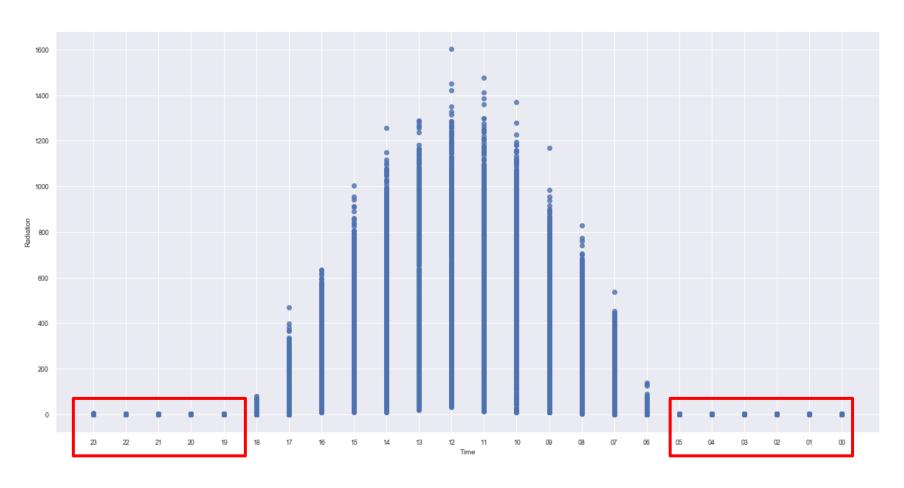
Time >> 시 / 분 / 초



상관관계 분석

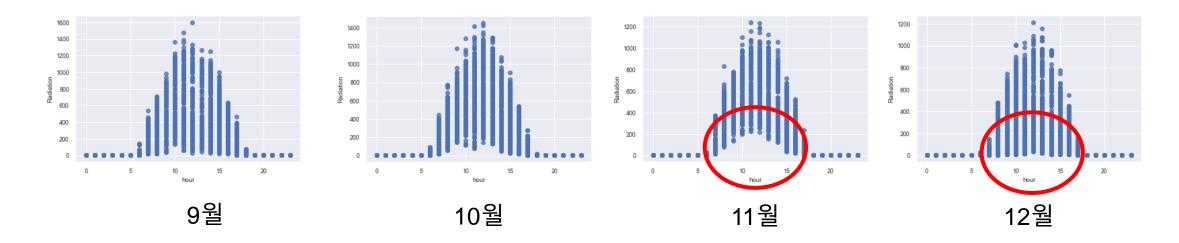
- Radiation 일사량과의 correlation
- 기온, 기압, 습도에서 유의미함을 찾음
- 시간과의 상관관계 값이 낮은 이유?





시간별 Scatter plot

월별 Scatter plot

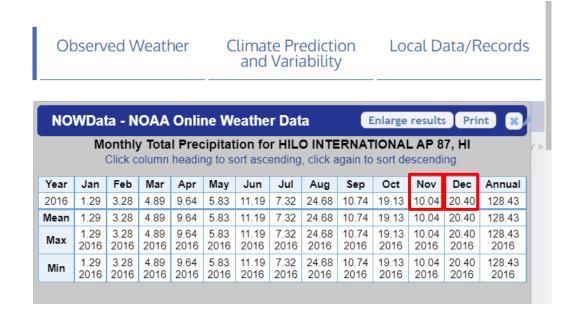


- 월별 Radiation 분포 양상의 차이
- 날씨와 관련이 있지 않을까?

Part 2 EDA



- 2016년 Hawaii 가까운 지역의 날씨 자료 검색
- 월별 평균 강수량 11월이 가장 낮고,
 12월이 가장 높음
- Column의 연관성 도출



Part 2 EDA

pyspark dataframe aggregation -평균

```
>>> avg temp.show()
                              >>> avg_hum.show()
                                                              >>> avg prs.show()
month
                 avg temp|
                                                avg hum|
                               |month|
                                                              month
                                                                                avg prs
    9| 53.68113792638533|
                                  12 | 79.5264576188143 |
                                                                 11| 30.44577981651425|
   10 | 52.46865434757964 |
                                   9|79.48577592018336|
                                                                 10|30.438462759324807|
   11| 50.78500724287784|
                                  10|78.94637796168234|
                                                                 9| 30.4320978832418|
   12|47.608892699657034|
                                                                 12|30.374427976481897|
                                   11 | 62.38495895702559 |
```

- 기온 - 습도 - 기압

Feature selection

- 기온(Temperature)

- 기압(Pressure)

- 습도(Humidity)

- 풍향(WindDirection)

- 풍속(Speed)

- month

- date

- year

hour

- minute

- second

- seminuter

- setminute

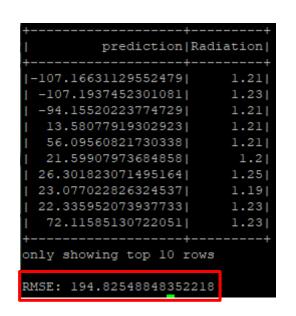
['Temperature', 'Pressure', 'Humidity', 'WindDirection(Degrees)', 'Speed', 'month', 'date', 'year', 'hour', 'minute', 'second', seminuter', 'setminute']



모델링

Pyspark Mllib

Linear regression



=> Random Forest

- RMSE 값 194.8254 >> 너무 높다 판단

전체 데이터 RF

```
prediction | Radiation |
 20.64930932403363|
                         1.21
                         1.23|
 20.64930932403363|
 21.88778013672256|
                         1.21
20.127617574793227|
                         1.21
20.127617574793227|
                         1.21|
|19.474924950520755|
                          1.2
|16.578980361202103|
                         1.25|
18.405581801487337
                         1.19|
18.595186123626906
                         1.23
20.342016086796868|
                         1.23
only showing top 10 rows
RMSE: 133.6635324751923
```

모델링

전체 시간대 RF

prediction|Radiation| 108.57711239279669 1.24 108.950429311068431 1.23 116.78631176427697| 1.21 135.4062482754064| 33.75 87.231 143.90313352002264| 446.2201864787821| 310.59 584.5928396749007 790.08 654.9495014654186| 804.28 836.2194301880897 855.63| 925.16 892.8198850886295 only showing top 10 rows RMSE: 181.61323818686026

```
prediction|Radiation|
 50.18108784538548|
                         1.21
  50.2170389099095|
                         1.21
                         1.29
53.272084988340886|
 42.50661521360623|
                         2.06
 67.57788723303783|
                         6.09
271.51620998675537|
                       177.01
 472.3500760493069|
                       548.94
 472.3500760493069|
                       365.99
 452.7766285037229|
                        374.21
 446.01171203165586
                       302.82
only showing top 10 rows
RMSE: 154.93211819140254
```

```
prediction | Radiation
  43.34546074136546|
                           1.2
                          1.18
    76.200744626504|
                          1.21
 104.69037924246558|
                          1.18
  57.04298694911089|
  47.77967353084034|
                          6.44
  193.23483129886391
                        191.05
                        282.43
  501.7442419006167|
  492.36735728227461
                        245.73
  502.39608969732891
                        307.57
 505.627692354623831
                        346.94
only showing top 10 rows
RMSE: 131.6274520058872
```

prediction|Radiation 77.84291007379369| 1.24 77.842910073793691 1.2 77.842910073793691 1.16 90.43177771392862| 8.49 94.679101091904951 15.96 187.822289841074021 552.33 427.55680404779986 637.19 427.556804047799861 692.01 447.272410206013431 261.18 432.74947262561926 448.57 only showing top 10 rows RMSE: 123.06854006562283

9월 10월 11월 12월

모델링

일사량 변동이 큰 시간대 별 RF

prediction|Radiation| 17.123040942678088| 1.21 17.123040942678088 1.23 20.414371164523956 1.21 1.21 16.924508947835257 36.20976165713893| 1.21 19.53429027003343| 1.2 1.25 20.414371164523956 20.414371164523956 1.19 17.123040942678088 1.23 1.23 19.314791002542698| only showing top 10 rows RMSE: 137.67213475551068

prediction|Radiation 36.11335242199591| 1.26 36.11335242199591 1.26 36.11335242199591| 1.29 36.11335242199591| 1.19 1.23 36.11335242199591 29.473709241058145| 1.27 36.11335242199591| 1.29 1.23 32.44610080825871| 32.44610080825871| 1.24 46.36307627320209| 1.22 only showing top 10 rows RMSE: 129.1776117055599

prediction|Radiation| 52.443078347447496| 1.21 1.22 52.4430783474474961 1.2 52.443078347447496| 52.443078347447496| 1.25 2.01 93.98133678110266| 239.77554076487422| 94.54 284.3854599459856| 262.97 286.57 473.70890316053741 382.93337223613065| 245.73 362.481 498.75376993111803| only showing top 10 rows RMSE: 101.16916858635068

prediction|Radiation| 45.778202954050826| 1.19 1.21 45.778202954050826| 1.21 45.778202954050826 45.778202954050826| 1.22 54.30203397472891| 1.2 32.62225948387426| 1.23 42.22907956072924| 1.22 41.93576960035581| 1.22 1.23 28.13418013174632| 44.20572866056205| 1.24 only showing top 10 rows RMSE: 93.07351381684322

9월 10월 11월

12월



Part 4 결과 분석

데이터 분석 결과

- 시간대별, 월별 데이터 분포가 다르게 나타남
- 영향을 주는 column의 특성, 일사량과의 연관성 발견

모델링 결과 분석

- · Random Forest 사용시 정확도 높아짐
- 시간대별, 월별 일사량 예측 결과 전체 데이터로 예측한 결과보다 RMSE 감소



개선 방향

- 기온이 높은 월의 데이터 모델 정확도 향상 필요



