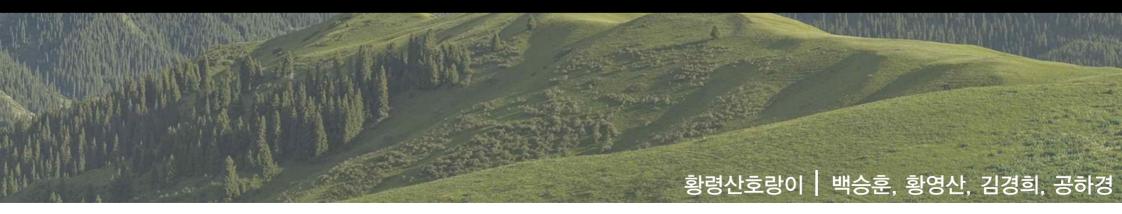


2021 날씨 빅데이터 콘테스트 24, 48시간 후 산사태 발생 예측 모델 개발



CONTENTS

01 공모배경

05 데이터 전처리(2)

98 데이터

06 분석기법 및 결과

03 데이터 전처리(1)

07 활용방안 및 기대효과

04 탐색적 자료 분석



산사태 발생 현황



세계 곳곳 이상기후 확산…'기록적인' 폭우·폭염에 시름 여지훈 기자 2021-07-18

경향신문

극한 넘어선 '이상기후', 더 세지고 잦아져…남 일일까 김한솔 기자 2021-07-19

세계일보

온난화가 불러온 역 대급 '물폭탄'... 산사 태 등 사고 속출 김유나, 남혜정 기자 2020-08-05

이상기후로 인한

산사태 발생 위험성이 높아지는 추세

산사태 피해 및 복구비 현황

15년도 이후 산사태 발생 빈도 점차 증가, 특히 2020년의 경우 매우 급증 따라서 산사태로 인한 피해복구비용 또한 증가 ↑



3 🎧 날씨 빅데이터 콘테스트

기존 산사태 예측 모델의 한계

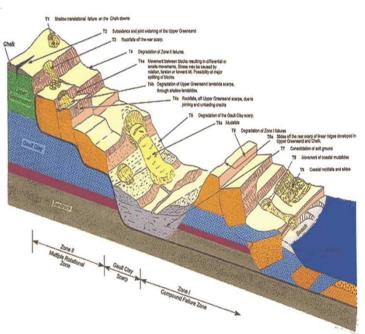
극심한 기후변화



중 · 단기 예측 정보 제공 불가



다각적인 요인들

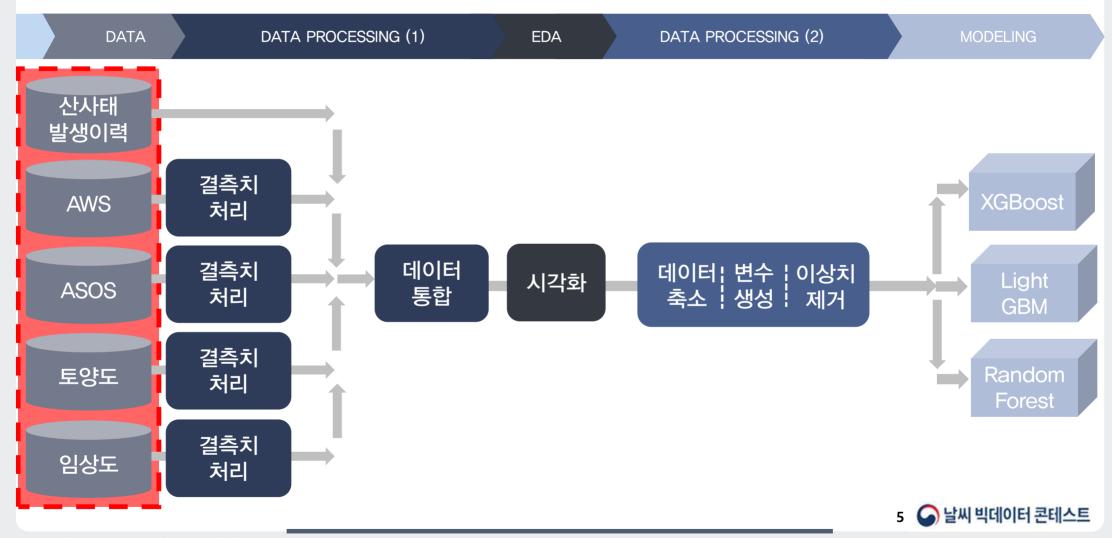


기상 데이터를 활용해 산사태 발생을 예측함으로써,

산사태로 인한 인명 및 재산피해 최소화 하는 것이 목표



전체 프로세스



제공 데이터 [산사태 발생이력, AWS, ASOS, 토양도, 임상도]

Variable	발생이력	Variable	AWS, ASOS	Variable	토양도, 임상도
tma	날짜	aws_min_ps	최저 해면기압	PRRCK_LARGE_CD	모암대코드
sd	시도	aws_avg_rhm	평균 상대습도	PRRCK_MDDL_CD	모암중코드: 화강암류, 반암류, 규장암류 등을 포함한 17개 의 모암중코드로 구분
sgg	시군구	aws_min_rhm	최저 상대습도	LOCTN_ALTTD	입지표고 : 입지토양도의 표고수치
umd	읍면동	aws_avg_ws	평균 풍속	LOCTN_GRDNT	입지경사도 : 입지토양도의 경사도
sum_cnt	산사태 발생 횟수 합	aws_max_ws	최대 풍속	CLZN_CD	기후대코드 : 온대북부, 온대중부, 온대남부, 난대
aws_sum_rn	합계 강수량	aws_max_ins_ws	최대 순간 풍속	TRGRP_TPCD	지형구분코드 : 계곡, 구릉지, 능선 등을 포함한 13개의 지형 구분코드로 구분
aws_sum_rn_dur	합계 강수 계속시간	asos_sum_rn	합계 강수량	SLDPT_TPCD	토심구분코드: 10, 20, 30
aws_hr1_max_rn	1시간 최다 강수량	asos_avg_rhm	평균 상대 습도	FRTP_CD	임상코드: 무립목지/비산림, 침엽수림, 활엽수림, 혼효림
aws_avg_tca	평균 전운량	asos_avg_td	평균 이슬점 온도	KOFTR_GROU_CD	수종그룹코드
aws_avg_ps	평균 해면기압	asos_avg_ws	평균 풍속	LDMARK_STNDA_CD	지형지물표준코드
aws_max_ps	최고 해면기압	SLTP_CD	토양형코드	SCSTX_CD	토성코드

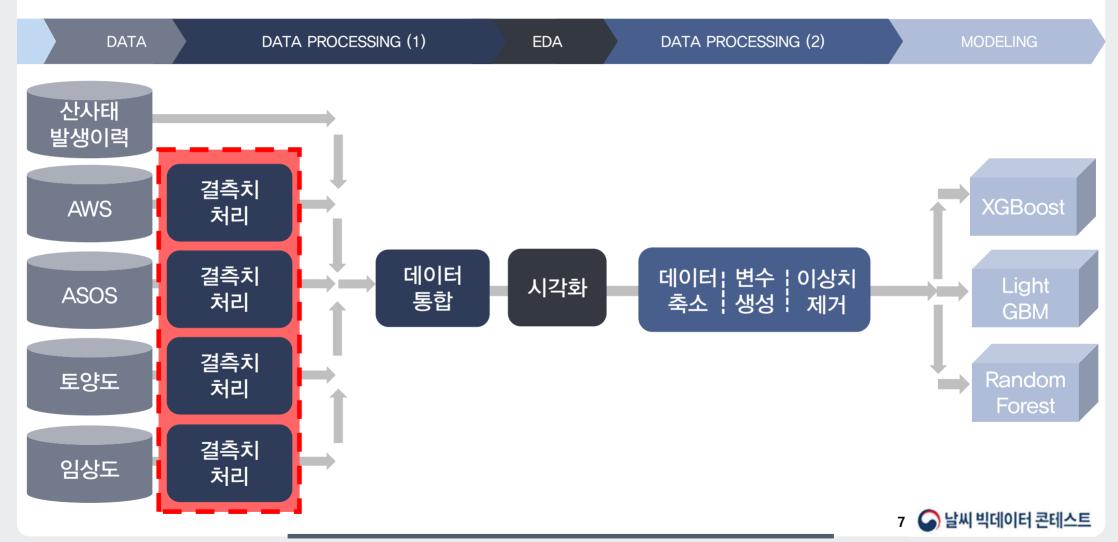
산사태 발생이력: 2011-2019년도 데이터

AWS/ASOS 기상 데이터: 637개의 읍,면,동에 대한 2011-2020년도 데이터

토양도, 임상도: 여러 폴리곤에 대한 정보를 개별 읍,면,동으로 대치하기 위해 수치형 자료는 평균 처리, 범주형 자료는 최빈값 처리를 하였음



전체 프로세스



결측치 처리 [AWS, ASOS, 토양도, 임상도]



AWS 관측지점 위치 (좌측)

ASOS 관측지점 위치 (우측)

토양도, 임상도

토양도, 임상도 공간데이터(.shp)와 경상도의 읍,면,동 공간데이터를 QGIS 응용프로그램을 통해 매칭 읍,면,동 폴리곤 내 토양도, 임상도 정보들을 mode 및 mean 처리를 하여 각 읍,면,동을 대표하는 값으로 변환 정보가 전혀 없는(결측) 지역은 QGIS를 통해 근방 지역의 값으로 대체

AWS, ASOS

QGIS를 이용하여 각 읍,면,동을 근처 기상관측지점으로 매칭 AWS, ASOS는 관측지점 위치 차이 존재 강수량 결측값(NA)는 0으로 처리 강수량 이외의 AWS, ASOS 기상데이터 변수들은 시계열 보간법 중 스플라인 보간법 사용

토양도, 임상도 shp 파일 에시 (좌측) 읍,면,동 (우측)

결측치 처리 [Train data(2011-2019), Test data(검증셋)]

통합 데이터

2011 ... 2019 2020 2021

1. 통합데이터 결측 처리

결측 시점: (2018-09-05, 2018-09-06, 2018-11-08)

(2020년-6월)

AWS 지점(822)에서 관측 값이 없는 날짜 존재. 따라서 해당지점의 66개 행 제거(매칭된 읍.면.동 22 곳)

2. 검증셋 결측 처리

2020년 6월 행정구역 코드에 존재하지 않는 읍,면,동 존재 (삼거동, 어곡동, 주진동, 시동, 덕곡동, 조와동)

결측 지점 보간법

AWS, ASOS

근처 지점번호로 대체

토양도, 임상도

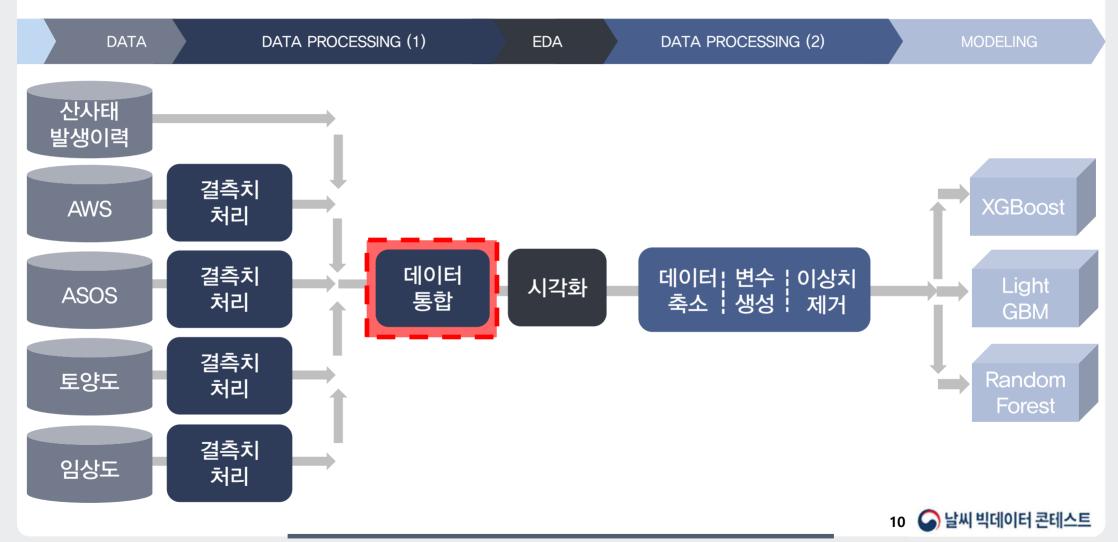
우측 표와 같이 QGIS로 얻은

근방 지역의 값으로 대체

읍면동	위도	경도	지점번호	경도(degree)	위도(degree)	노장해발고도(m)	지점명(한글)	예보구역코드	법정동코드
거제시 삼거동	34.84016937	128.6498519	294	128.60453	34.88818	45.4	거제	11H20403	4831025021
양산시 어곡동	35.39546521	129.0078563	257	129.02	35.3072	14.9	양산시	11H20102	4833031027
양산시 주진동	35.39554054	129.1414906	257	129.02	35.3072	14.9	양산시	11H20102	4833031027
경주시 시동	35.75767944	129.271876	283	129.2009092	35.81746056	37.6	경주시	11H10202	4713011300
김천시 덕곡동	36.11937334	128.1518407	279	128.32	36.13	48.8	구미	11H10602	4719010800
영주시 조와동	36.86708332	128.6328298	272	128.51696	36.87188	210.8	영주	11H10401	4721025021



전체 프로세스



데이터 통합 [산사태 발생이력, AWS, ASOS, 토양도, 임상도]

1) AWS, ASOS 데이터(2011-2019): 변수 12, 4개 2) 토양도, 임상도 데이터: 변수 9, 3개 3) 산사태 발생이력 데이터: 변수 1개

		-					
tma	stn_id	sum_rn	•••	max_ws	avg_ws	시도명칭	Y
2011-01-01	115	28		19	11.1	경상북도	포
2011-01-02	115	7		8.5	6	경상북도	포
2011-01-03	115	16		10.2	6.5	경상북도	포
2011-01-04	115	0		6.8	3.1	경상북도	포
2011-01-05	115	1		6.9	4.1	경상북도	포
2011-01-06	115	16.5		11.5	5.9	경상북도	포
2011-01-07	115	7		6.7	2.4	경상북도	포
2011-01-08	115	0		7.6	3.3	경상북도	포
2011-01-09	115	12.5		14	5.4	경상북도	포

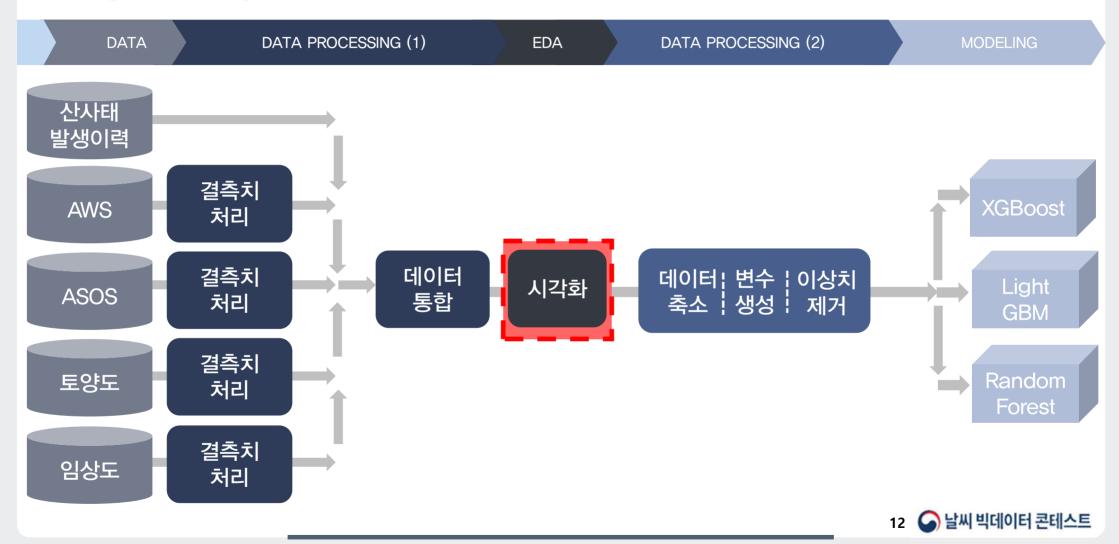
시도명칭	시군구명칭	읍면동명칭	mode_PRRC K_LARG	mode_PRRC K_MDDL	mode_MAP_L ABEL
경상북도	포항시 남구	구룡포읍	2	24	NPD-13C
경상북도	포항시 남구	연일읍	2	22	NMM-24C
경상북도	포항시 남구	오천읍	1	12	NMM-24C
경상북도	포항시 남구	대송면	1	12	NMM-24C
경상북도	포항시 남구	동해면	2	25	NEB-24C
경상북도	포항시 남구	장기면	2	24	NMM-24C
경상북도	포항시 남구	호미곶면	2	24	NEB-24C
경상북도	포항시 남구	송도동	2	23	NPT-25C
경상북도	포항시 남구	청림동	2	24	NMM-24C

tma	sd	sgg	umd	sum_cnt	sum_hpa		
20110709	경상남도	밀양시	내일동	1	1.2		
20110709	경상남도	밀양시	단장면	4	3.7		
20110709	경상남도	밀양시	무안면	5	4.2		
20110709	경상남도	밀양시	부북면	6	7.8		
20110709	경상남도	밀양시	산외면	1	2		
20110709	경상남도	밀양시	상동면	6	12.5		
20110709	경상남도	밀양시	청도면	2	1.3		
20110709	경상남도	밀양시	초동면	1	4		
20110709	경상남도	사천시	곤명면	1	1		

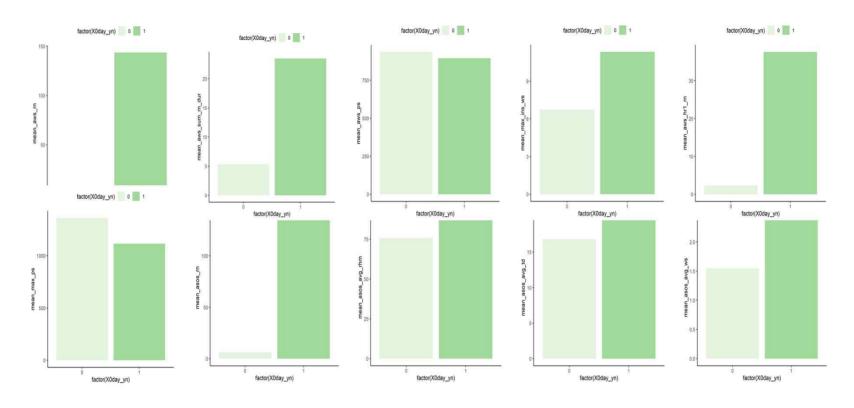
기간: 2011-01-01 ~ 2019-12-31: 제공 데이터에서 변수 총 29개 사용

date	sd	sgg	umd	code	stn_id.aws.	stn_id.asos.	aws_sum _rn	aws_sum_r n_dur	aws_hr1_m ax_rn	aws_avg_t ca	 mode_SLDP T_TPCD	mode_FR TP_CD	mode_KOFT R_GROU	mode_LDMAR K_STN
2011-01-01	경상북도	포항시 남구	구룡포읍	3701111	138	138	3.5	20	1	0	 20	2	30	J12230
2011-01-01	경상북도	포항시 남구	연일읍	3701112	138	138	3.5	20	1	0	 20	1	77	J12377
2011-01-01	경상북도	포항시 남구	오천읍	3701113	138	138	3.5	20	1	0	 10	2	30	J12230
2011-01-01	경상북도	포항시 남구	대송면	3701131	138	138	3.5	20	1	0	 10	2	30	J12230
2011-01-01	경상북도	포항시 남구	동해면	3701132	138	138	3.5	20	1	0	 20	2	30	J12230
2011-01-01	경상북도	포항시 남구	장기면	3701133	138	138	3.5	20	1	0	 10	2	30	J12230

전체 프로세스



시각화 [AWS, ASOS]

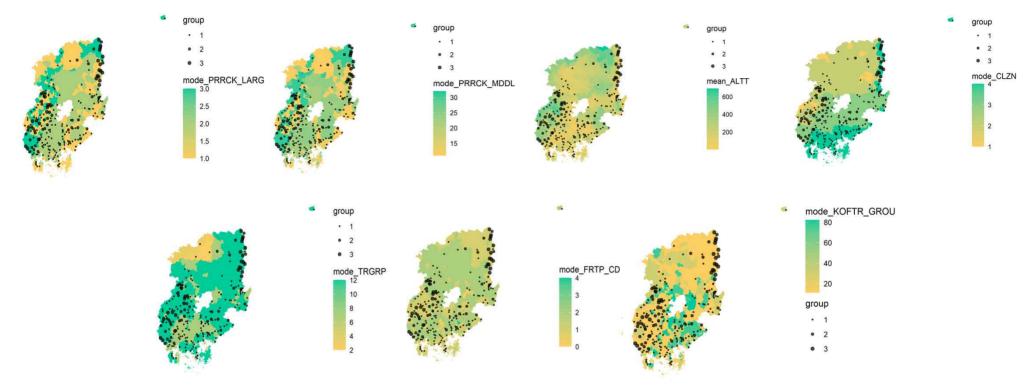


AWS, ASOS 기상데이터의 변수에 대하여 산사태 발생(0), 산사태 발생하지 않음(1)으로 그래프를 그려본 결과,

대다수의 날씨변수들은 산사태가 발생한 날과 발생하지 않은 날 간의 차이가 매우 크게 두드러짐



시각화 [토양도,임상도]



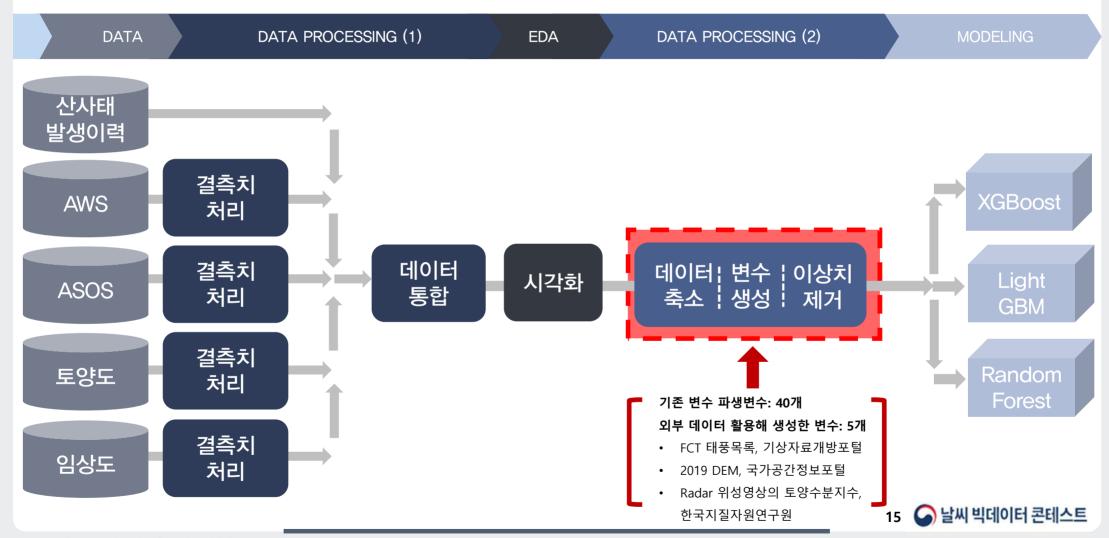
경상도 지형에 토양도, 임상도 변수들에 따라 색의 차이를 주고 산사태 발생이력 데이터의 sum_cnt 빈도에 따라 점의 크기에 변화를 주었음.

위 시각화 지도 그래프를 통해 7pg에 해당하는 토양도, 임상도 변수의 값들이 각 지역에 고루 분포되어 있음을 확인.

단, 7pg에 서술되지 않은 기타 토양도, 임상도 변수들은 단일 값이거나 분산이 매우 적어 분석과정에 포함하지 않았음.



전체 프로세스



데이터 축소

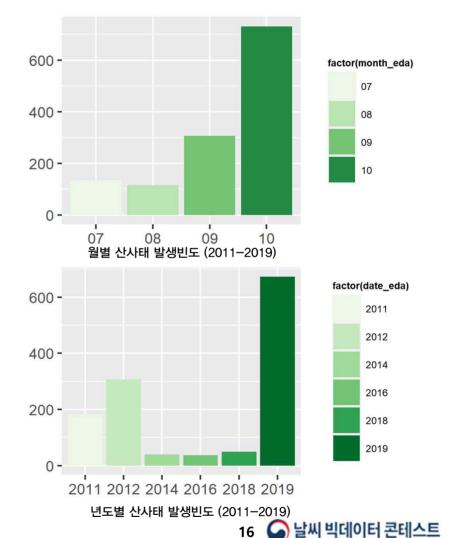
2,093,753 rows

Date	 SLTP
2011-01-01	 2
2011-01-02	 2
2011-01-03	 2
2011-01-04	 2
2019-12-24	 82
2019-12-25	 82
2019-12-26	 82
2019-12-27	 82
2019-12-28	 82
2019-12-29	 16
2019-12-30	 16
2019-12-31	 16

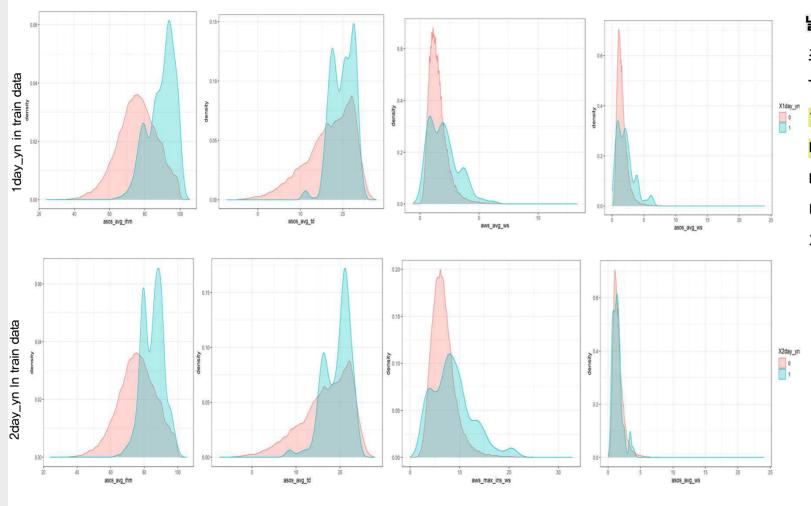
우측 그래프에 따르면 2013년, 2015년, 2017년도
산사태가 한번도 발생하지 않았음을 알 수 있음.
따라서 통합 데이터에서 2013, 2015, 2017년도를 제외하였으며 지난 10년간 6월 ~10월에만 산사태가 발생하였음에 기초하여 6월~10월 데이터만 사용.

583,358 rows

Date		SLTP
2011-06-01		2
2011-06-02	•••	2
2011-06-03		2
2019-10-28		16
2019-10-29		16
2019-10-30		16
2019-10-31		16



diff변수, rn_rn_dur 변수 [AWS, ASOS]



날씨 파생변수: diff 변수 (32개)

좌측 그래프에 따르면

Train 데이터에서 기상자료 변수들의

1day_yn 와 2day_yn 양상이 확연히

다름을 볼 수 있음.

따라서 모든 기상자료 변수들(16개)에

대하여 1일 후 와 2일 후 차이를

계산한 diff변수를 각각 생성

1일 뒤 diff 변수명: Diff1 - Diff16

2일 뒤 diff 변수명: Diff2_1 - Diff2_16

날씨 파생변수: rn_rn_dur 변수

rn_rn_dur = sum_rn x sum_rn_dur

(누적강우량 * 강수 지속기간)

의미: 강수량에 지속 시간으로

가중치를 준 변수를 추가 생성

17 🕝 날씨 빅데이터 콘테스트

태풍 변수 [특보데이터]

FCT 태풍발생목록

태풍명	영문명	영향도	발생	소멸
에어리	AERE	없음	2011-05-07	2011-05-12
송다	SONGDA	없	2011-05-22	2011-05-29
사리카	SARIKA	아 전	2011-06-10	2011-06-11
하이마	HAIMA	암	2011-06-21	2011-06-25
메아리	MEARI	접 징 경 영	2011-06-22	2011-06-27
망온	MA-ON	암	2011-07-12	2011-07-24
도카게	TOKAGE	암	2011-07-15	2011-07-16
고구마	KOGUMA	아 잡	2021-06-12	2021-06-13
참피	СНАМРІ	없음	2021-06-23	2021-06-27

산사태 발생이력

Date	umd	Sum_cnt
20110709	내일동	1
20110709	단장면	4
20110709	무안면	5
20110709	부북면	6
20110709	산외면	1
20110709	상동면	6
20110709	청도면	2
20110709	초동면	1
20110709	곤명면	1

특보데이터 변수: storm 변수 (3개)

2011-2019년도 태풍발생목록과

2011-2019년도 산사태 발생이력 데이터를

병합한 결과 7일 이상 지속된 태풍의

태풍 소멸일 근방에서 경상도 지역의 산사태가

발생했음을 알 수 있었음.

→ stormday, stormday1, stormday2 으로

각각 태풍 소멸일, 태풍소멸 1일 전, 태풍소멸 2일 전에

이진 분류(0: 발생하지 않음,1: 발생 했음)값을 채운

변수를 생성하였음.

추가 데이터(1) :

기상자료개방포털 〉데이터 〉기상예보 〉태풍예보

토양군, 지형 습윤, 토양수분 변수 [임상도, DEM, Radar 위성영상]

임상도 파생변수: 토양군 soil_group

국내 산림토양 분류방식 (SLTP_CD 토양형 코드 기반)

01-06 갈색산림토양(B)

07-09 적/황색산림토양(R_Y)

10-14 암적색산림토양(DR)

15-16 회갈색산림토양(GrB)

...

지형 습윤지수: humidity

 $TW I = \ln(\frac{A}{\tan(B)})$

A: 상류지역 수분기여면적(*m*²) B: 경사도(radian)

지형 습윤

(TWI; Topographic Wetness Index)

강우가 어떤 방향으로 흘러가는지, 흘러갔을 때의 양을 수치적으로 표현한 변수

'0' 에 가까울수록 강우시 집수되는 물의 양이 적음 커질수록 강우시 집수되는 물의 양이 많음 ←→

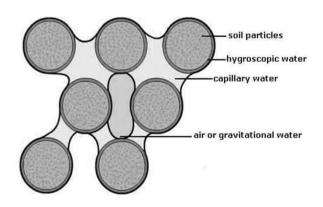
> TWI의 값이 작을수록 건조한 토양환경 값이 클수록 습윤한 토양환경

추가 데이터(2): 2019 DEM 국가공간정보포털 토양 수분 변수: soil_moisture

토양 수분은 경사면 내의 압력을 변화

시켜 경사면을 불안정하게 하기 때문에 산사태를 유발할 수 있음 . 결과적으로, 무거운 물을 함유한 토양 및 암석은 중 력에 더 큰 영향을 받음. 따라서 과도한 토양수분은 산사태의

가장 주요한 원인으로 작용



추가데이터(3): Radar 위성영상의 토양수분지수

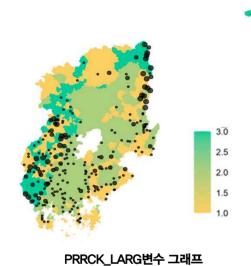
한국지질자원연구원

☑ 날씨 빅데이터 콘테스트



diff_LARG변수 [토양도]

토양도 파생변수[1]: diff_LARG변수



시군구명칭	읍면동명칭	읍면동코드	diff_LARG
거창군	위천면	3839035	0
경산시	남부동	3710054	1
경산시	남산면	3710034	0
구미시	원평1동	3705052	2
구미시	원평2동	3705053	0

	화성암(1)	퇴적암(2)	변성암(3)
화 성 암 (1)	0	1	2
퇴 적 암 (2)	1	0	1
변 성 암 (3)	2	1	0

diff LARG 변수 형태

diff_LARG 변수 생성 규칙

PRRCK_LARG변수(모암대 코드)를 위 그래프로 표현한 결과, 모암대 코드가 바뀌는 경계 지역에서 특히 많은 산사태가 발생함을 확인.

따라서 근방 읍,면,동에 대해 모암대 코드의 변화에 따라 범주형 변수인 diff_LARG 생성하였음.

diff_LARG 생성 규칙은 상단 우측 표과 같음. (PRRCK LARG 코드 값: 1(화성암), 2(퇴적암), 3(변성암))

공극 관련 변수 [토양도]

토성분류	유효 저류능[2]	최소 침투율[3]	침투율[4]	NRCS 분류[5]	유효 공극률[6]
SCSTX_CD	Storage_capacity	Min_infiltration	inflitration	NRCS	Porosity
SL(Sandy loam)(01)	0.25	1.02	5.715	В	0.412
L(Loam)(02)	0.19	0.52	5.715	В	0.434
Si(Silty clay)(08)	0.09	0.04	0.635	D	0.423
LS(Loamy sand)(09)	0.31	2.41	9.525	А	0.401
S(Sand)(10)	0.35	8.27	9.525	A	0.417
C(Clay)(11)	0.08	0.02	0.635	D	0.385

※ 공극: 암석 또는 토양 입자 사이의 틈

토양의 <mark>공극</mark>에

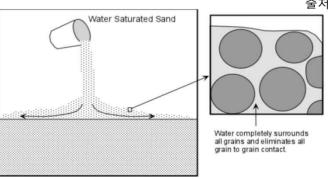
물이 충분히 들어가게 되면

입자 간의 미끌어짐 발생

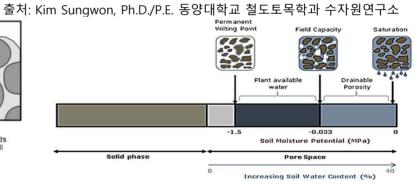
Ex. 모래 (참고 - 우측 그림[1])

약한 습윤상태: 서로 잘 붙어있음

강한 습윤상태: 입자들이 사방으로 흩어짐

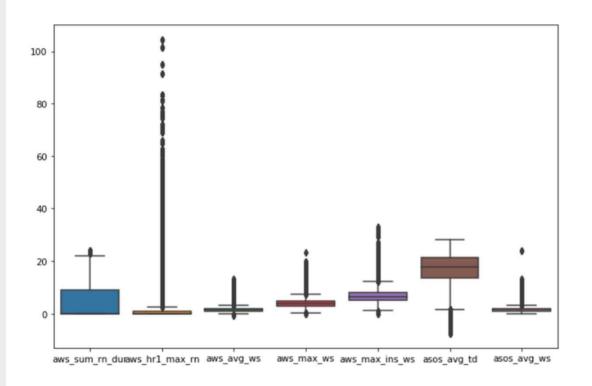


그림[1] 모래 습윤상태에 따른 입자구조도



그림[2] 보편적 토양 습윤상태에 따른 입자구조도

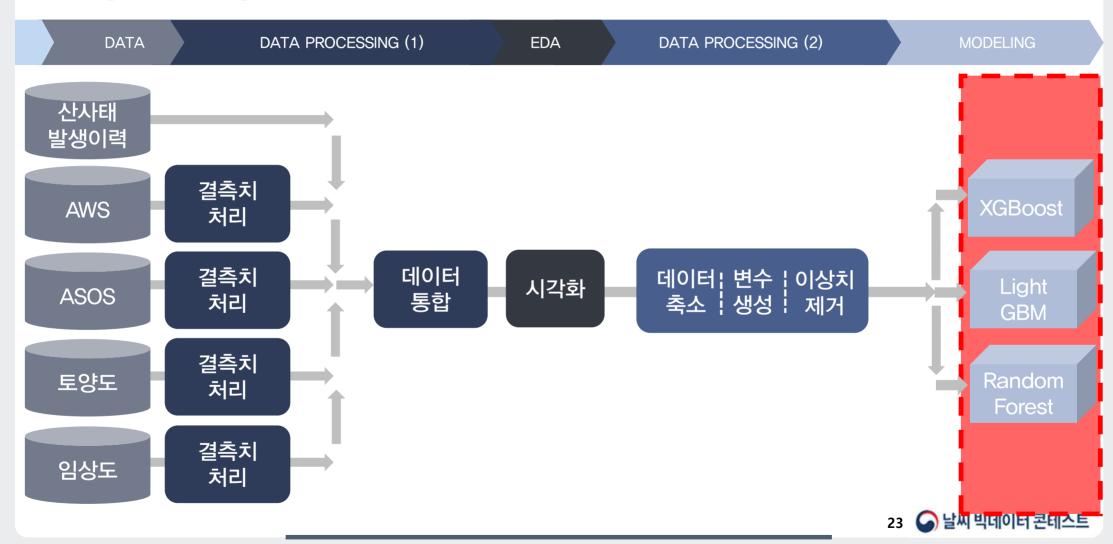
이상치 제거



- 이상치 데이터(Outlier)는 모델의 성능을 떨어뜨 리는 불필요한 요소
- 연속형 변수 64개에서 25개 이상의 변수에서 outlier가 존재하는 행을 제거
- 총 583358개의 행 중에서 3947개의 행(전체 의 0.7%)이 제거됨

〈연속형 변수의 boxplot〉

전체 프로세스



모델 선택

데이터의 문제

1. 데이터 불균형

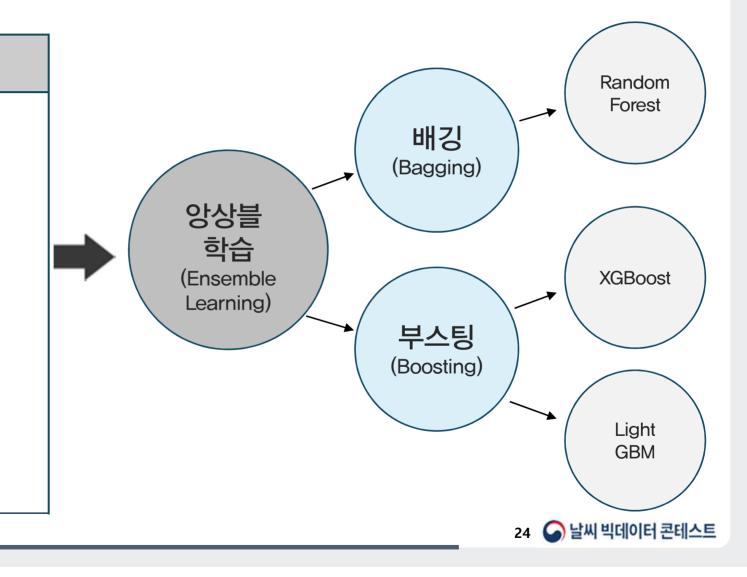
산사태 비율: 0.06% 심한 unbalanced 데이터 셋

2. 일반화 성능

Target=1 의 수가 매우 작아 오버피팅의 가능성이 매우 높음

3. 속도와 성능

산사태 발생 특성상 빠른 데이터 업데이트와 모델의 성능이 필요



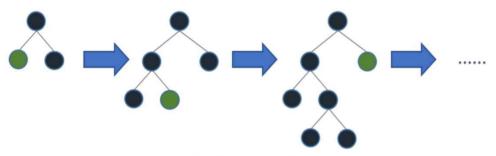
모델 별 성능 비교

Model	UnderSa	mpling O	UnderSampling X		
	Accuracy	모델 종합 정확도 (CSI)	Accuracy	모델 종합 정확도 (CSI)	
Random Forest	93.13%	16.79%	93.39%	18.09%	
LightGBM	93.65%	18.34%	95.1%	21.58%	
XGBoost	94.59%	19.04%	94.32%	19.23%	

참가번호 210124 의 모델 평가 점수는 Accuracy 95.1%, CSI는 21.58% 입니다.

최종 선정 모델

LightGBM



Leaf-wise tree growth

- 산사태 데이터는 불균형 데이터이기 때문에 undersampling을 시행해봄
- 이 때 cross validation(K = 5)을 이용하였고, 각 모델은 파라 미터를 조정해가며 가장 높은 성능을 기준으로 함
- 결과적으로 Undersampling을 하지 않는 LightGBM의 성능이 가장 좋음
- cut-off는 각각 0.0017(24시간), 0.0018(48시간)

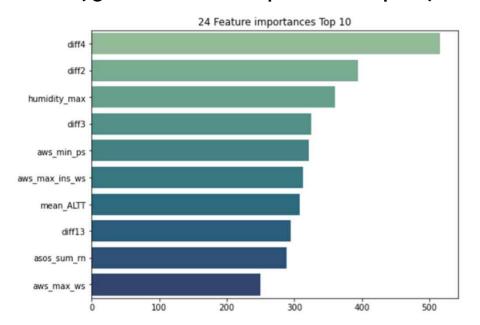
Parameter	24시간	48시간
n_estimators	500	500
learning_rate	0.01	0.01
max_depth	7	7
num_leaves	48	32
subsample	0.65	0.8
colsample_bytree	0.79	1

[Hyperparameters of LGBM]

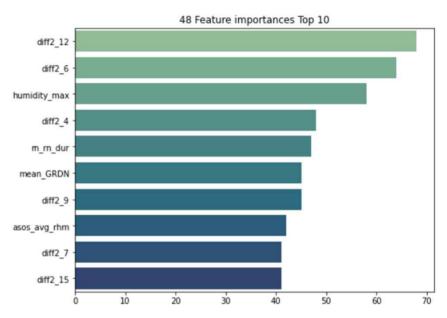


모델 해석

⟨lgbm24 Feature Importance Top 10⟩



⟨lgbm48 Feature Importance Top 10⟩



- 최종적으로 선택된 모델의 중요변수로는 평균 전운량의 차이, 합계 강수 계속시간의 차이, 최대 순간 풍속의 차 이, 최고 해면기압의 차이, 지형습윤지수 등이 있음
- 주로 1일 뒤, 2일 뒤 기상 데이터와의 차이 변수인 diff변수가 중요 변수로 보여짐

산사태 피해 예방을 위한 예측 경보 시스템







산사태예측정보제공

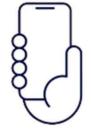


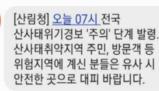
















기대 효과



기후 변화에 따른 산사태 취약성 보완



단기(24시간 전), 중기(48시간 전)에 산사태 예측정보를 제공함으로 선제적인 주민 대피를 통해 인명피해 최소화



실시간 산사태 경보 알림을 통해 피해 금액, 피해 면적 최소화

참고문헌

[참고문헌]

- 산사태 정의, 산사태정보시스템, (2021년 06월 20일), https://sansatai.forest.go.kr/forecast/introMap.do.
- 김기흥, 정혜련, 박재현, 마호섭, (2011), 경남지역 산사태 발생지의 강우 및 지형특성분석, 한국환경복원기술학회지 14(2), 33-45,
- 김석우, 전근우, 김진학, 김민식, 김민석, (2012), 2011년 집중호우로 인한 산사태 발생특성분석, 한국임학회지 101(1), 28-35.
- 강경희. (2019). Random Forest 모델을 이용한 경상북도 상주 지역과 강원도 진부 지역의 산사태 취약성 분석 (국내석사학위논문). 세종대학교 대학원. 서울.
- 이지훈. (2021년 07월 18일), 세계 곳곳 이상기후 확산 '기록적인 '폭우. 폭염에 시름. 데일리뉴스. https://www.idailynews.co.kr/news/view.php?idx=83241
- 김한솔. (2021년 07월 19일). 극한 넘어선 '이상기후' 더 세지고 잦아져…남 일일까. 경향신문, https://www.khan.co.kr/world/world-general/article/202107192134005
- 김유나, 남혜정, (2020년 08월 05일), 온난화가 불러온 역대급 '물폭탄' ··· 산사태 등 사고 속출, 세계일보, http://m.seqve.com/view/20200804523391
- 김용준. (2020년 10월 06일). [지난 3년 여름의 경고] ② 산사태 비 오는 지역은 무조건 산사태 특보?...천재인가. 인재인가. KBS.

https://n.news.naver.com/article/056/0010911210

[추가 데이터]

- 국도교통부, 행정구역 읍면동(법정동). 2021년 6월 30일 검색. 국가공간정보포털
- 기상청, 태풍목록, 2021년 7월 20일 검색, 기상자료개방포털
- 국도교통부 국토지리정보원, 수치표고모델(DEM) 90M. 2021년 8월 1일 검색. 국가공간정보포털
- 한국지질자원연구원, Radar 위성영상의 토양수분지수(2019). 2021 8월 1일 검색. Envbigdata

Q&A

감사합니다