

BÁO CÁO TIẾN ĐỘ ĐỒ ÁN MÔN HỌC

THÔNG TIN NHÓM

STT	MSSV	Họ và tên	Vai trò
1	19522065	Nguyễn Thị Minh Phương	Nhóm trưởng
2	19521882	Chu Hà Thảo Ngân	
3	19522397	Thái Minh Triết	

TÊN ĐỀ TÀI: XÂY DỰNG MÔ HÌNH DỰ ĐOÁN TÁC ĐỘNG CỦA CÁC YẾU TỐ TỰ NHIÊN ĐẾN TÌNH TRẠNG SẠT LỞ ĐẤT.

I. Nguồn dữ liệu

1. Bộ dữ liệu chính

- Tên bộ dữ liệu: *Global Landslide Catalog (GLC)*.
- Nguồn dữ liệu: <https://data.nasa.gov/Earth-Science/Global-Landslide-Catalog-Export/dd9e-wu2v>
- Kích thước: 11033 x 31, bao gồm 22 categorical features và 9 numerical features.
- Mô tả bộ dữ liệu: bộ dữ liệu *GLC* được biên soạn năm 2007 tại *NASA Goddard Space Flight Center*. Bộ dữ liệu chứa thông tin liên quan đến các sự kiện sạt lở đất trên toàn thế giới khoảng từ năm 1988 đến năm 2017.

2. Dữ liệu thu thập thêm

2.1. Thời tiết

- Nguồn dữ liệu: <https://www.visualcrossing.com/weather-api>
- Kích thước: 11033 x 19, bao gồm 2 categorical features và 17 numerical features.

- Mô tả nguồn dữ liệu: dữ liệu được thu thập thêm từ Weather API của *VisualCrossing*. Bộ dữ liệu có chứa các thông tin liên quan đến thời tiết tại thời điểm và địa điểm xảy ra sự kiện sạt lở đất trong bộ *GLC*.

2.2. Elevation

- Nguồn dữ liệu: <https://developers.airmap.com/docs/elevation-api>

- Kích thước: 11033 x 1, là numerical feature.

- Mô tả nguồn dữ liệu: dữ liệu được thu thập từ Elevation API của trang web Airmap. Chứa thông tin về độ cao so với mực nước biển cho hầu hết các vị trí địa lý trên Trái Đất. Độ đo được sử dụng là mét. Độ phân giải không gian: 1 arc-second (khoảng 30 mét).

2.3. Continent

- Nguồn dữ liệu: *package pycountry_convert*.

- Kích thước: 11033 x 1, là categorical feature.

- Mô tả nguồn dữ liệu: package sử dụng dữ liệu từ Wikipedia, cho phép thực hiện việc chuyển đổi giữa tên quốc gia (chuẩn ISO) sang mã quốc gia và châu lục. Các giá trị của châu lục thu được từ bộ dữ liệu gồm có:

- **Asia:** Châu Á
- **North America:** Bắc Mỹ
- **South America:** Nam Mỹ
- **Europe:** Châu Âu
- **Africa:** Châu Phi
- **Oceania:** Châu Đại Dương

2.4. Season

- Nguồn tham khảo: <https://www.nationalgeographic.org/encyclopedia/season/>

- Kích thước: 11033 x 1, là categorical feature.

- Mô tả nguồn dữ liệu: dữ liệu về mùa tại sự kiện xảy ra sạt lở đất được thu thập dựa vào thời điểm và vị trí địa lý nó so với đường xích đạo.

Vị trí	Thời gian	Season
Bắc bán cầu	21/03 – 20/06	Spring
	21/06 – 22/09	Summer
	23/09 – 20/12	Autumn
	21/12 – 20/03 năm sau	Winter
Nam bán cầu	21/03 – 20/06	Autumn
	21/06 – 22/09	Winter
	23/09 – 20/12	Spring
	21/12 – 20/03 năm sau	Summer

2.5. Mật độ dân số

- Nguồn thu thập: Bộ dữ liệu **Gridded Population of the World Version 4.11**

- <https://sedac.ciesin.columbia.edu/data/set/gpw-v4-population-density-rev11>
- Google Earth Engine: https://developers.google.com/earth-engine/datasets/catalog/CIESIN_GPWv411_GPW_UNWPP-Adjusted_Population_Density

- Kích thước: 11033 x 5, là numerical feature.

- Mô tả nguồn dữ liệu: dữ liệu về mật độ dân số (người/km²) của thế giới trong các năm 2000, 2005, 2010, 2015 và 2020 được thu thập và ước tính bởi *Center for International Earth Science Information Network, Columbia University* và *NASA Socioeconomic Data and Applications Center*, dữ liệu được hiệu chỉnh để phù hợp với dự báo triển vọng dân số thế giới năm 2015 (The World Population Prospects: 2015 Revision) của Liên Hợp Quốc.

- Cách thức thu thập: dữ liệu được thu thập sử dụng thư viện *ee* trong package *geemap* trong python để truy xuất đến bộ dữ liệu **GPWv411:UN-Adjusted Population Density** trong Google Earth Engine, sau đó lấy dữ liệu về mật độ dân số ở từng năm theo kinh độ và vĩ độ.

2.6. Độ che phủ rừng

- Nguồn thu thập: bộ dữ liệu **Hansen Global Forest Change v1.8 (2000-2020)**

- <https://data.globalforestwatch.org/documents/134f92e59f344549947a3eade9d80783/explore>

- Google Earth Engine: https://developers.google.com/earth-engine/datasets/catalog/UMD_hansen_global_forest_change_2020_v1_8

- Kích thước: 11033 x 3, là numerical feature.

- Mô tả nguồn dữ liệu: dữ liệu về độ che phủ rừng toàn cầu năm 2000 và sự thay đổi (tăng và giảm) về độ che phủ rừng từ năm 2000 đến năm 2020. Dữ liệu được tổng hợp và số hóa từ hình ảnh vệ tinh Landsat 7.

- Cách thức thu thập: dữ liệu được thu thập sử dụng thư viện *ee* trong package *geemap* trong python để truy xuất đến bộ dữ liệu **Hansen Global Forest Change v1.8 (2000-2020)** trong Google Earth Engine, sau đó lấy dữ liệu về độ che phủ rừng theo kinh độ và vĩ độ.

2.7. Kết cấu của đất (Soil Texture)

- Nguồn thu thập: Bộ dữ liệu **OpenLandMap Soil Texture Class (USDA System)**

- Google Earth Engine: https://developers.google.com/earth-engine/datasets/catalog/OpenLandMap_SOL_SOL_TEXTURE-CLASS_USDA-TT_M_v02

- Kích thước: 11033 x 6, là categorical feature.

- Mô tả nguồn dữ liệu: dữ liệu về 12 loại kết cấu đất (soil texture) theo hệ thống USDA ở 6 độ sâu: 0 cm, 10 cm, 30 cm, 60 cm, 100 cm và 200 cm.

- Cách thức thu thập: dữ liệu được thu thập sử dụng thư viện *ee* trong package *geemap* trong python để truy xuất đến bộ dữ liệu **OpenLandMap Soil Texture Class (USDA System)** trong Google Earth Engine, sau đó lấy dữ liệu về loại kết cấu đất ở từng độ sâu theo kinh độ và vĩ độ.

3. Bộ dữ liệu kết hợp cuối cùng

- Bộ dữ liệu cuối cùng được kết hợp từ các bộ dữ liệu thành phần trên, dựa trên khóa là vị trí và thời điểm xảy ra sự kiện sạt lở đất (“latitude”, “longitude”, “event_date”).

- Kích thước toàn bộ của bộ dữ liệu: 11033 x 67, bao gồm 32 categorical features và 35 numerical features.

II. Thông tin mô tả bộ dữ liệu

1. Global Landslide Catalog (GLC)

Index	Name	Data-type	Description	Data range
1	source_name	object	Tên báo đưa tin	
2	source_link	object	Liên kết dẫn đến tin	
3	event_id	int64	Mã sự kiện sạt lở đất	
4	event_date	object	Giờ/ngày/tháng/năm diễn ra sạt lở đất	1988-11-07 to 2017-09-28
5	event_time	float64	Giờ diễn ra sạt lở đất	
6	event_title	object	Tiêu đề tin tức sạt lở đất	
7	location_description	object	Mô tả thông tin vị trí sạt lở	
8	location_accuracy	object	Khoảng cách chênh lệch giữa vị trí ghi nhận so với vị trí thực tế	<ul style="list-style-type: none"> ▪ unknown (542) ▪ 5km (3,178) ▪ 10km (1,435) ▪ 25km (1,470) ▪ exact (1,386) ▪ 1km (2,185) ▪ 50km (794) ▪ 250km (16) ▪ 100km (25)
9	event_description	object	Mô tả sự kiện sạt lở đất	<ul style="list-style-type: none"> ▪ landslide (7,648) ▪ mudslide (2,100) ▪ complex (232) ▪ rock_fall (671) ▪ debris_flow (194) ▪ riverbank_collapse (37) ▪ other (68) ▪ unknown (38) ▪ lahar (7) ▪ snow_avalanche (15) ▪ creep (5) ▪ earth_flow (7) ▪ translational_slide (9) ▪ topple (1)

10	landslide_category	object	Loại sạt lở đất	<ul style="list-style-type: none"> ▪ landslide (7,648) ▪ mudslide (2,100) ▪ complex (232) ▪ rock_fall (671) ▪ debris_flow (194) ▪ riverbank_collaps e (37) ▪ other (68) ▪ unknown (38) ▪ lahar (7) ▪ snow_avalanche (15) ▪ creep (5) ▪ earth_flow (7) ▪ translational_slide (9) ▪ topple (1)
11	landslide_trigger	object	Nguyên nhân gây ra sạt lở đất	<ul style="list-style-type: none"> ▪ downpour (4671) ▪ rain (2556) ▪ unknown (1689) ▪ continuous_rain (746) ▪ tropical_cyclone (561) ▪ snowfall_snowmel t (132) ▪ monsoon (129) ▪ mining (93) ▪ earthquake (89) ▪ construction (79) ▪ flooding (72) ▪ no_apparent_trig ger (43) ▪ freeze_thaw (41) ▪ other (24) ▪ dam_embankmen t_collapse (11) ▪ leaking_pipe (10) ▪ vibration (1) ▪ volcano (1)

12	landslide_size	object	Mức độ sạt lở đất	<ul style="list-style-type: none"> ▪ large (750) ▪ small (2,767) ▪ medium (6,551) ▪ unknown (851) ▪ very_large (102) ▪ catastrophic (3)
13	landslide_setting	object	Môi trường xung quanh vị trí sạt lở đất	<ul style="list-style-type: none"> ▪ mine (157) ▪ unknown (6,291) ▪ above_road (3,104) ▪ urban (264) ▪ natural_slope (531) ▪ engineered_slope (22) ▪ below_road (199) ▪ above_river (149) ▪ retaining_wall (48) ▪ other (50) ▪ above_coast (20) ▪ bluff (48) ▪ burned_area (28) ▪ deforested_slope (53)
14	fatality_count	float64	Số lượng người tử vong	0 to 5000
15	injury_count	float64	Số lượng người thương vong	0 to 374
16	storm_name	object	Tên cơn bão xảy ra trước khi sạt lở	
17	photo_link	object	Đường dẫn tới hình ảnh khu vực bị sạt lở	
18	notes	object	Ghi chú	
19	event_import_source	object	Nguồn cung cấp sự kiện sạt lở	<ul style="list-style-type: none"> ▪ glc (9,379) ▪ test (90)

				▪ Included...found somewhere else (1)
20	event_import_id	float64	Mã cung cấp sự kiện sạt lở	
21	country_name	object	Tên quốc gia nơi xảy ra sự kiện	
22	country_code	object	Mã quốc gia	
23	admin_division_name	object	Tên đơn vị hành chính	
24	admin_division_population	float64	Dân số của đơn vị hành chính	0 to 13M
25	gazeteer_closest_point	object	Vị trí trên bản đồ gần nơi xảy ra sạt lở nhất	
26	gazeteer_distance	float64	Khoảng cách từ "gazeteer_closest_point" tới nơi xảy ra sạt lở	3e-5 to 215.45 (km)
27	submitted_date	object	Ngày nộp/hoàn thành sample trên dataset	2014-04-01 to 2017-11-21
28	created_date	object	Ngày tạo sample trên dataset	2017-11-20 to 2017-12-20
29	last_edited_date	object	Ngày cuối cùng chỉnh sửa sample trên dataset	2018-02-15
30	latitude	float64	Vĩ độ nơi xảy ra sạt lở	-46.77 to 72.62
31	longitude	float64	Kinh độ nơi xảy ra sạt lở	-179.98 to 179.99

2. Weather

Index	Name	Datatype	Description	Data range
32	tempmax	float64	Nhiệt độ cao nhất trong ngày tại địa điểm sạt lở	-22.1 to 45.1 (°C)
33	tempmin	float64	Nhiệt độ thấp nhất trong ngày tại địa điểm sạt lở	-27.1 to 32.1 (°C)

34	temp	float64	Nhiệt độ tại địa điểm sạt lở	-24.8 to 38.4 (°C)
35	feelslikemax	float64	Nhiệt độ cảm thấy cao nhất	-22.1 to 51.9 (°C)
36	feelslikemin	float64	Nhiệt độ cảm thấy thấp nhất	-32.3 to 41.0 (°C)
37	feelslike	float64	Nhiệt độ cảm thấy	-28.0 to 43.6 (°C)
38	dew	float64	Điểm sương (nhiệt độ hóa sương)	-27.9 to 175.5 (°C)
39	humidity	float64	Độ ẩm tương đối	14.69 to 100.0 (%)
40	precip	float64	Lượng mưa tại thời điểm “datetime”	0.0 to 498.0 (mm)
41	precipcover	float64	Tỉ lệ số giờ trong ngày có lượng mưa khác 0	0.0 to 100.0 (%)
42	windgust	float64	Gió giật	0.0 to 316.8 (m/s)
43	windspeed	float64	Tốc độ gió trung bình trên 1 phút	0.0 to 222.1 (m/s)
44	winddir	float64	Hướng gió	0.0 to 360.0 (m/s)
45	pressure	float64	Áp suất khí quyển theo mực nước biển hoặc khí áp (tính theo đơn vị millibars hay hectopascals)	951.8 to 1084.0 (Pa)
46	cloudcover	float64	Tỉ lệ bầu trời bị che phủ bởi mây	0.0 to 100.0 (%)
47	visibility	float64	Tầm nhìn xa	0.0 to 1030.1 1.0 (km)
48	moonphase	float64	Tỉ lệ hình dạng của mặt trăng theo chu kỳ	0.0 to 1.0
49	conditions	object	Điều kiện thời tiết	▪ Partially cloudy (1,360)

				<ul style="list-style-type: none"> ▪ Rain, Partially cloudy (4,455) ▪ Rain, Overcast (2,356) ▪ Rain (527) ▪ Overcast (333) ▪ Clear (569) ▪ Snow, Partially cloudy (107) ▪ Rain, Fog (4) ▪ Snow (26) ▪ Snow, Overcast (39)
50	stations	object	Trạm thời tiết	

3. Các features khác

Index	Name	Datatype	Description	Data range
51	elevation	int64	Độ cao so với mực nước biển (mét)	-407.0 to 6916.0 (m)
52	continent	object	Châu lục	<ul style="list-style-type: none"> ▪ Asia (4,930) ▪ North America (4,454) ▪ South America (490) ▪ Africa (235) ▪ Europe (570) ▪ Oceania (354)
53	season	object	Mùa	<ul style="list-style-type: none"> ▪ summer (3,732) ▪ winter (2,522)

				<ul style="list-style-type: none"> ▪ autumn (2,358) ▪ spring (2,421)
54	treecover2000	int64	Phần trăm độ phủ tán cây tại khu vực trong năm 2000 (áp dụng với thảm thực vật cao trên 5 mét)	0 to 100 (%)
55	loss	int64	Sự thay đổi từ forest sang non-forest trong khoảng thời gian ghi nhận.	<ul style="list-style-type: none"> ▪ 0 ▪ 1
56	gain	int64	Sự thay đổi toàn bộ từ non-forest sang forest trong khoảng thời gian từ năm 2000 - 2012	<ul style="list-style-type: none"> ▪ 0 ▪ 1
57	soil_texture_0	object	Loại kết cấu đất ở bề mặt	<ul style="list-style-type: none"> ▪ Lo ▪ ClLo ▪ SaClLo ▪ SaLo ▪ SaCl ▪ SiLo ▪ Cl ▪ LoSa ▪ SiClLo ▪ Sa
58	soil_texture_10	object	Loại kết cấu đất ở độ sâu 10 cm	
59	soil_texture_30	object	Loại kết cấu đất ở độ sâu 30 cm	
60	soil_texture_60	object	Loại kết cấu đất ở độ sâu 60 cm	
61	soil_texture_100	object	Loại kết cấu đất ở độ sâu 100 cm	
62	soil_texture_200	object	Loại kết cấu đất ở độ sâu 200 cm	
63	population_density_2000	float64	Mật độ dân số tại khu vực vào năm 2000	0.0 to 110884.89 (người/km ²)
64	population_density_2005	float64	Mật độ dân số tại khu vực vào năm 2005	0.0 to 116474.90 (người/km ²)

65	population_density_2010	float64	Mật độ dân số tại khu vực vào năm 2010	0.0 to 121829.19 (người/km ²)
66	population_density_2015	float64	Mật độ dân số tại khu vực vào năm 2015	0.0 to 126877.76 (người/km ²)
67	population_density_2020	float64	Mật độ dân số tại khu vực vào năm 2020	0.0 to 130805.85 (người/km ²)

III. Exploratory Data Analysis

1. Thông kê mô tả bộ dữ liệu

	count	mean	std	min	25%	50%	75%	max
event_id	11033	5598.953141	3249.228647	1	2785	5563	8435	11221
event_time	0	NaN	NaN	NaN	NaN	NaN	NaN	NaN
fatality_count	9648	3.219424	59.886178	0	0	0	1	5000
injury_count	5359	0.751819	8.458955	0	0	0	0	374
event_import_id	9471	4798.56307	2789.125559	-111.1673	2386.5	4773	7189.5	9669
admin_division_population	9471	157760.0458	829734.5446	0	1963	7365	34021	12691840
gazeteer_distance	9471	11.873689	15.598228	0.00003	2.363845	6.25487	15.81561	215.4489
longitude	11033	2.520441	100.908393	-179.980766	-107.8717	19.6946	93.948	179.9914
latitude	11033	25.881887	20.415054	-46.7748	13.9176	30.5345	40.866259	72.6275
tempmax	11004	18.857534	10.49672	-22.1	11.2	21.1	27.9	45.1
tempmin	11004	12.979198	9.323941	-27.1	4.8	13.2	22.1	32.1
temp	9776	17.725378	8.279178	-24.8	10.8	19.5	25.1	38.4
feelslikemax	11001	19.957177	12.40151	-22.1	11.2	21.1	30	51.9
feelslikemin	11001	12.307445	10.235743	-32.3	2.3	13.2	22.1	41
feelslike	9773	17.842014	9.817092	-28	10.4	19.5	25.5	43.6
dew	9740	14.202977	8.887815	-27.9	7.3	14.9	22.6	175.5
humidity	9740	81.949669	13.122496	14.69	76.9175	85.02	90.95	100
precip	9776	30.272759	52.780018	0	0.3	10.4	35.9725	498
precipcover	11004	26.585559	30.337672	0	0	13.64	45.83	100
windgust	3826	50.90149	18.46986	0	38.025	50	59.4	316.8
windspeed	9752	20.844924	12.993564	0	12.6	18.4	27.6	222.1

winddir	9269	176.340533	64.694292	0	132.7	175.9	219.4	360
pressure	8312	1010.184312	7.630745	951.8	1005.9	1010.1	1014.3	1084
cloudcover	9776	66.812848	30.089074	0	48.6	76	91.5	100
visibility	9698	18.535482	74.701551	0	7.2	10.8	14.9	1030.1
moonphase	11033	0.505495	0.304086	0	0.26	0.5	0.76	1
elevation	11033	714.618599	844.603545	-407	83	347	1120	6916

2. Kiểm tra giá trị khuyết

Tên thuộc tính	Số lượng missing value
event_time	11033
notes	10716
storm_name	10456
photo_link	9537
windgust	7207
injury_count	5674
pressure	2721
winddir	1764
admin_division_name	1637
country_code	1564
event_import_source	1563
gazeteer_closest_point	1563
country_name	1562
event_import_id	1562
gazeteer_distance	1562
admin_division_population	1562
fatality_count	1385
visibility	1335

dew	1293
humidity	1293
windspeed	1281
feelslike	1260
precip	1257
temp	1257
conditions	1257
cloudcover	1257
stations	1228
event_description	862
source_link	846
soil_texture_60	345
soil_texture_10	345
soil_texture_0	345
soil_texture_200	345
soil_texture_100	345
soil_texture_30	345
location_description	102
landslide_setting	69
feelslikemin	32

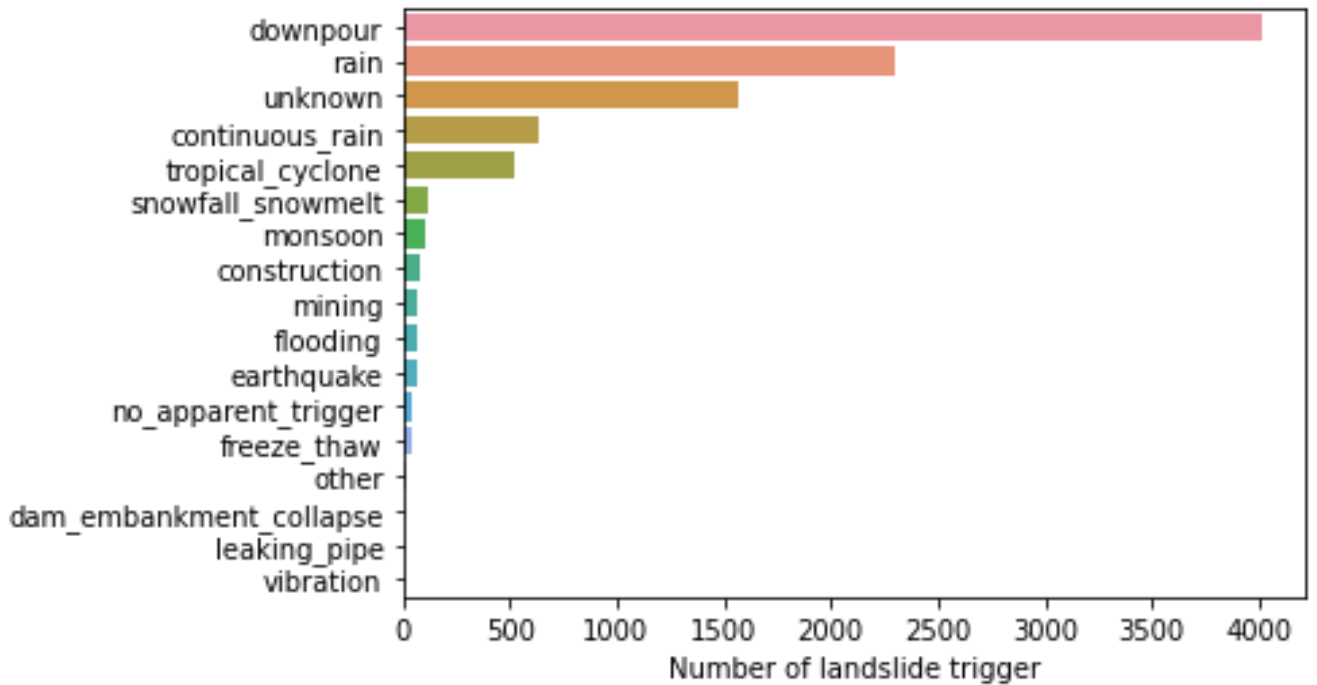
feelslikemax	32
tempmax	29
precipcover	29
tempmin	29
landslide_trigger	23
population_density_2005	21
population_density_2000	21
population_density_2010	21
population_density_2015	21
population_density_2020	21
loss	14
submitted_date	10
landslide_size	9

location_accuracy	2
created_date	1
landslide_category	1
moonphase	0
event_id	0
event_date	0
event_title	0
last_edited_date	0
latitude	0
longitude	0
gain	0
treecover2000	0
season	0
continent	0
elevation	0
source_name	0

Tỉ lệ dữ liệu bị khuyết giá trị: 12.18 %

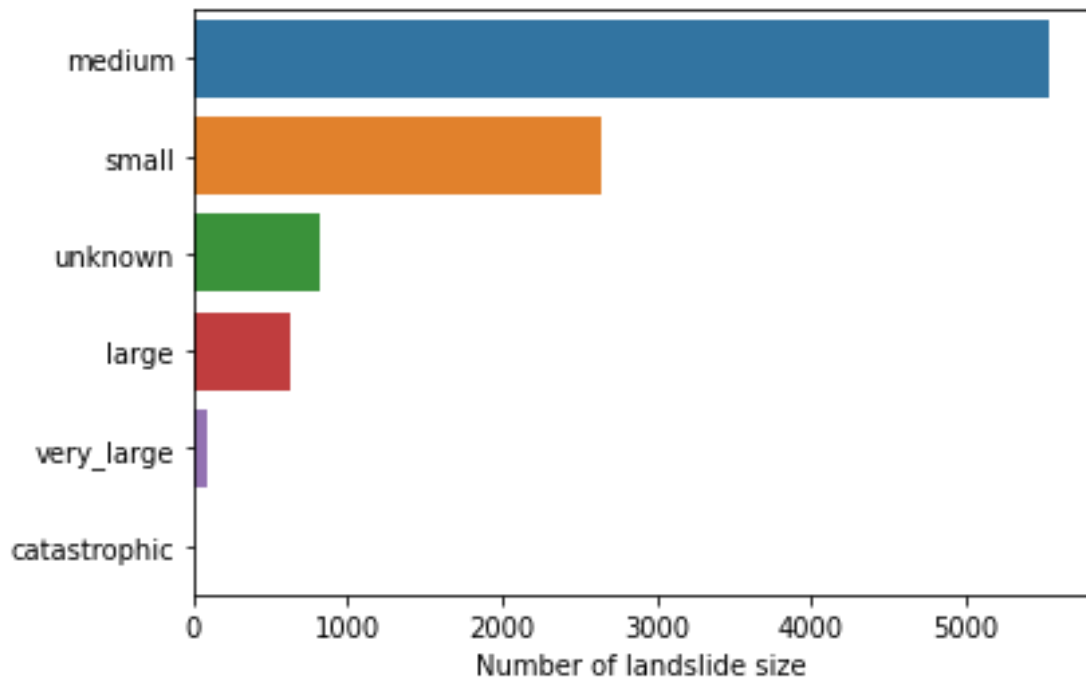
3. EDA dữ liệu *Global Landslide Catalog*

- Target: 'landslide_trigger' (nguyên nhân gây sạt lở)



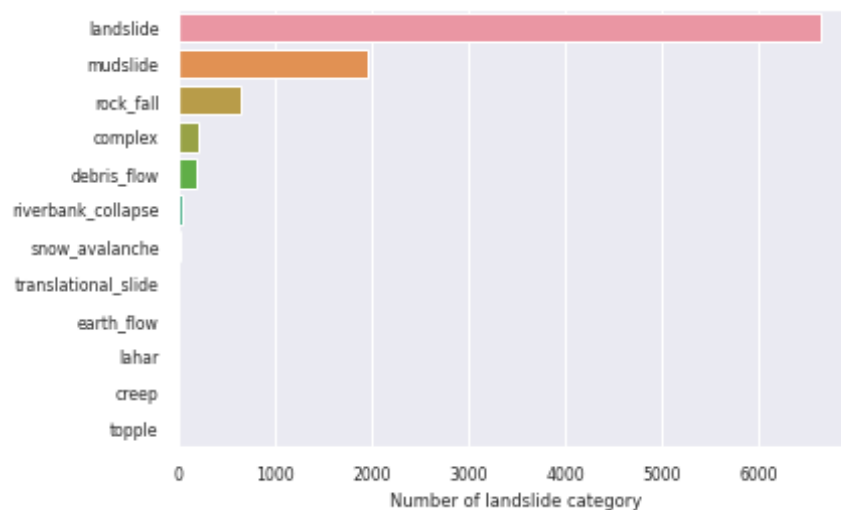
- Số lượng label ban đầu: 18
- Gom một số label về chung 1 label (ví dụ gom các loại mưa về chung label 'rain')
- Số lượng label sau khi xử lý: 12
- Số lượng label 'rain' mất cân bằng => sử dụng phương pháp giải quyết mất cân bằng như SMOTE, ADASYN

- **Target: 'landslide_size' (mức độ sạt lở)**

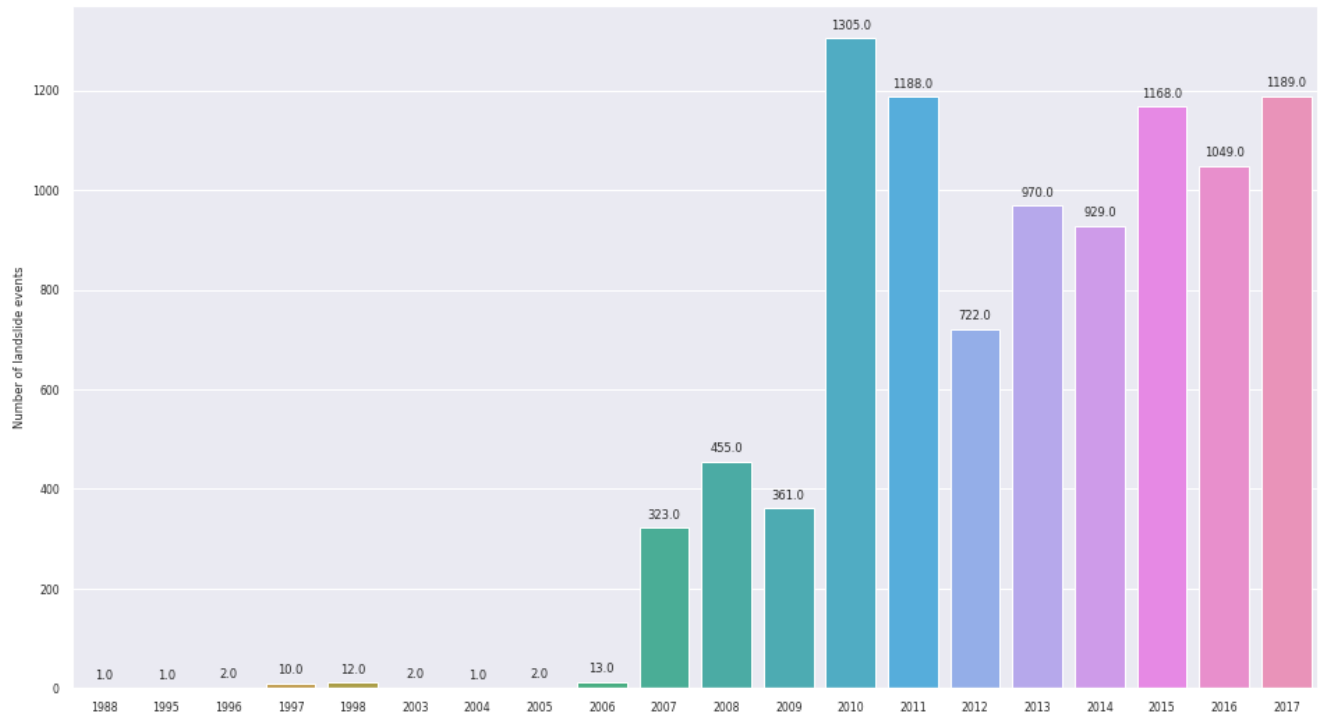


- Số lượng label ban đầu: 6
- Gom label ‘*very_large*’ chung với label ‘*large*’, bỏ label ‘*catastrophic*’
- Số lượng label sau khi xử lý: 4

- **Loại sạt lở đất (landslide_category)**



- **Số lượng những sự kiện sạt lở diễn ra hằng năm**



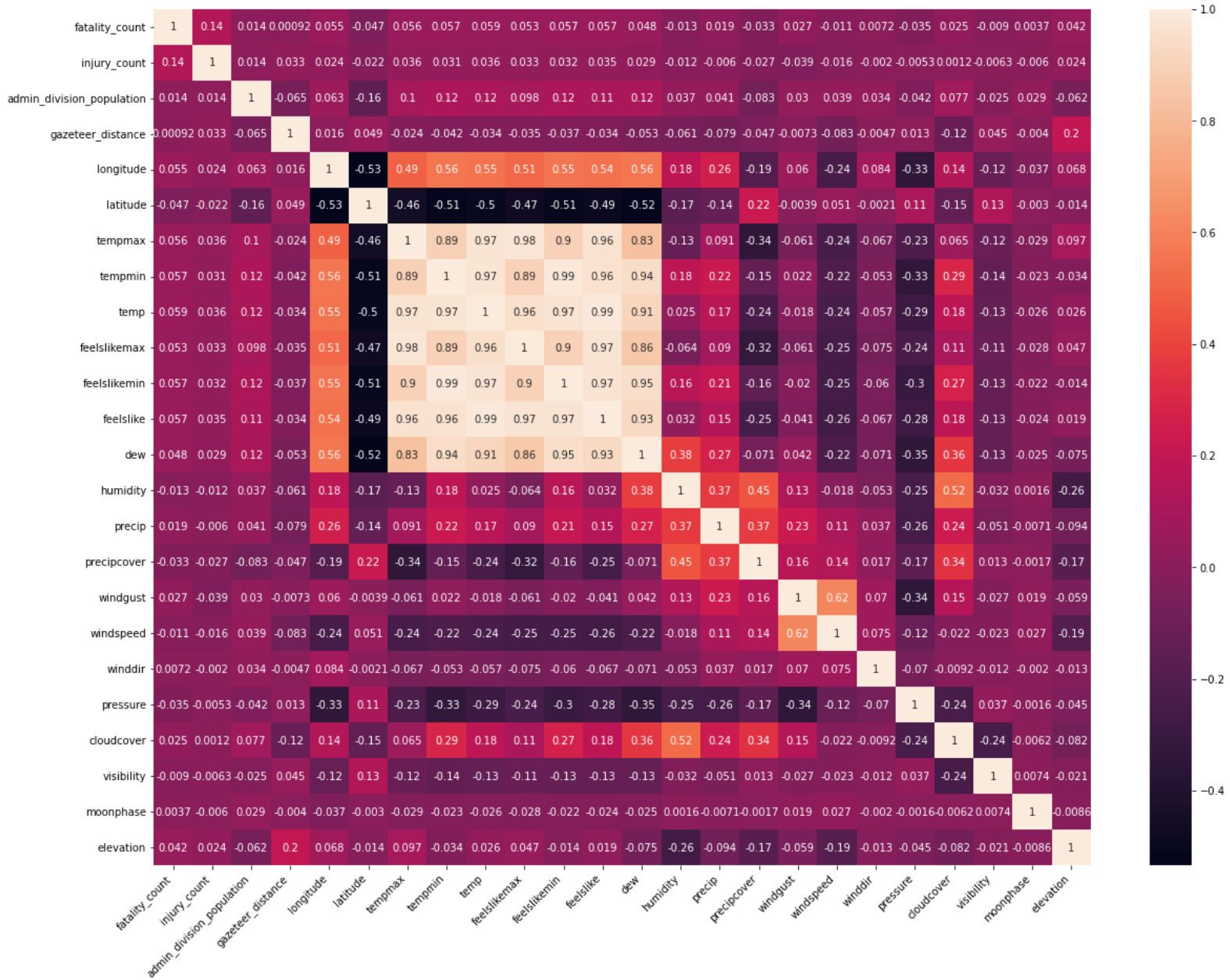
- Các sự kiện sạt lở đất trên toàn thế giới khoảng từ năm 1988 đến năm 2017
- NASA bắt đầu thu thập các thông tin liên quan sạt lở và biên soạn bộ dữ liệu từ năm 2007

- **Mã trận tương quan giữa thương-tử vong và dân số địa phương**

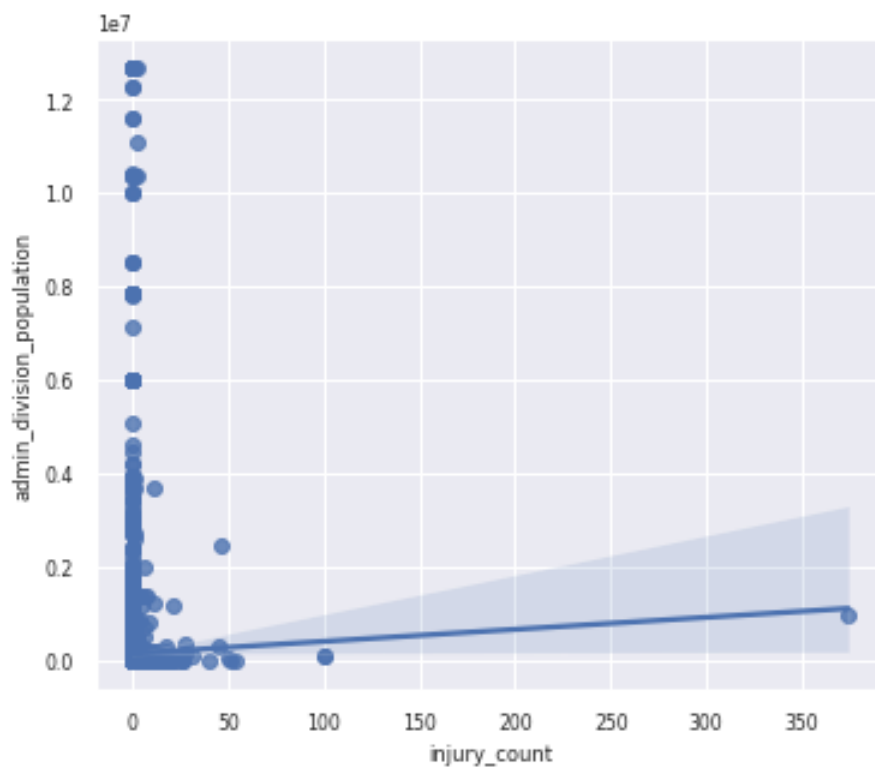
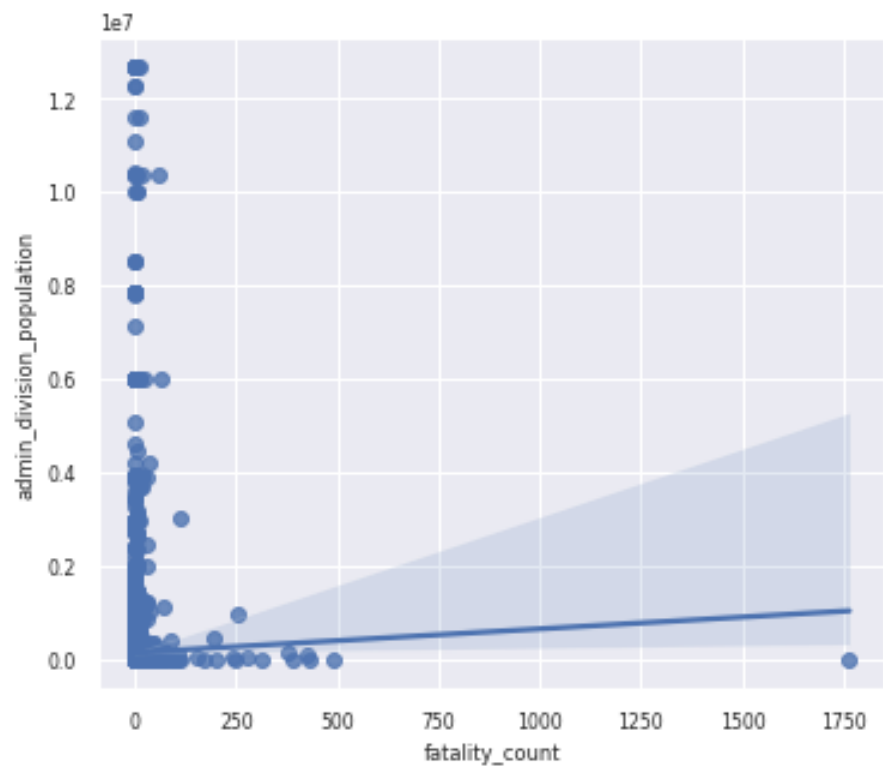
Giữa thương-tử vong và dân số địa phương (đơn vị hành chính) không có tương quan. Khu vực dân số tuy đông nhưng số thương-tử vong do sạt lở cũng thấp.

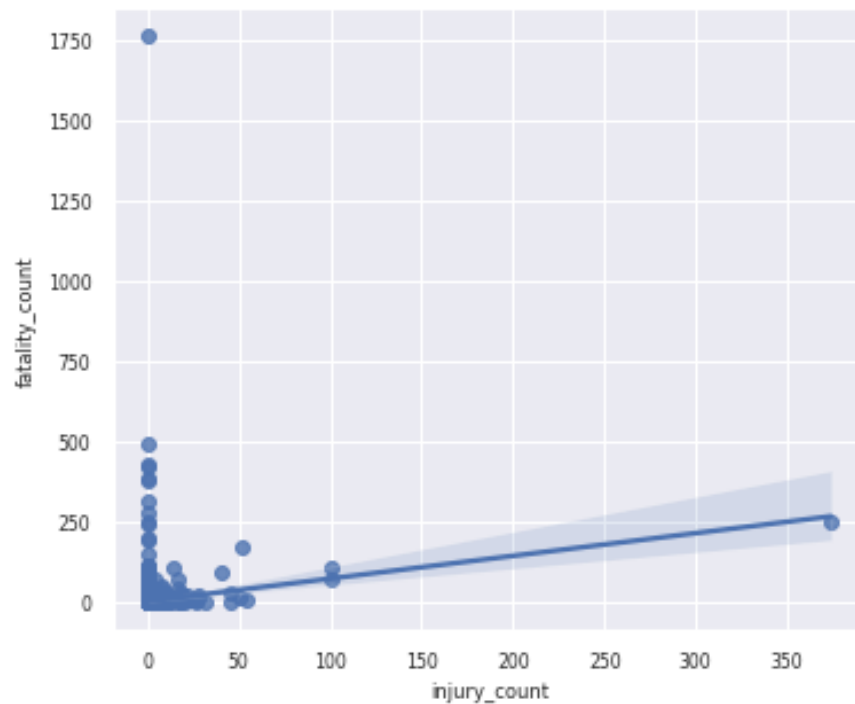


- Ma trận tương quan các thuộc tính

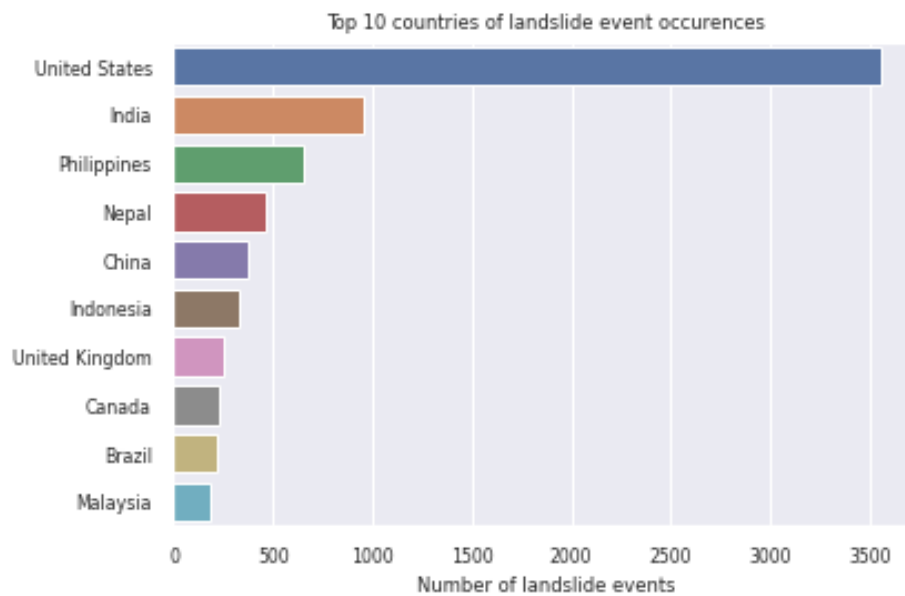


- Regplot thương-tử vong và dân số địa phương





- **Top 10 quốc gia xảy ra sạt lở nhiều nhất**

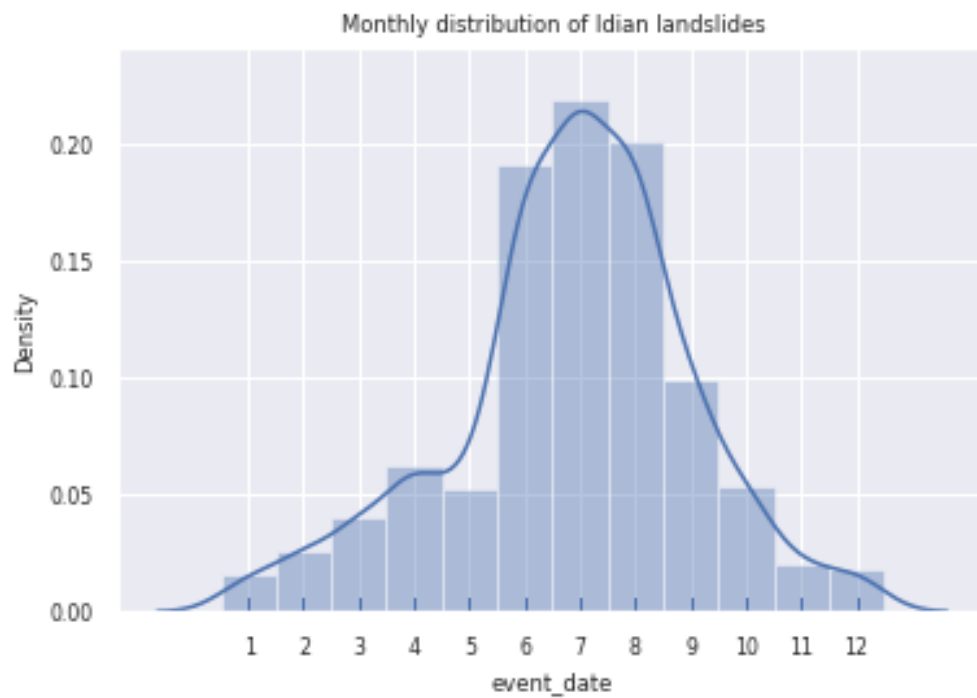
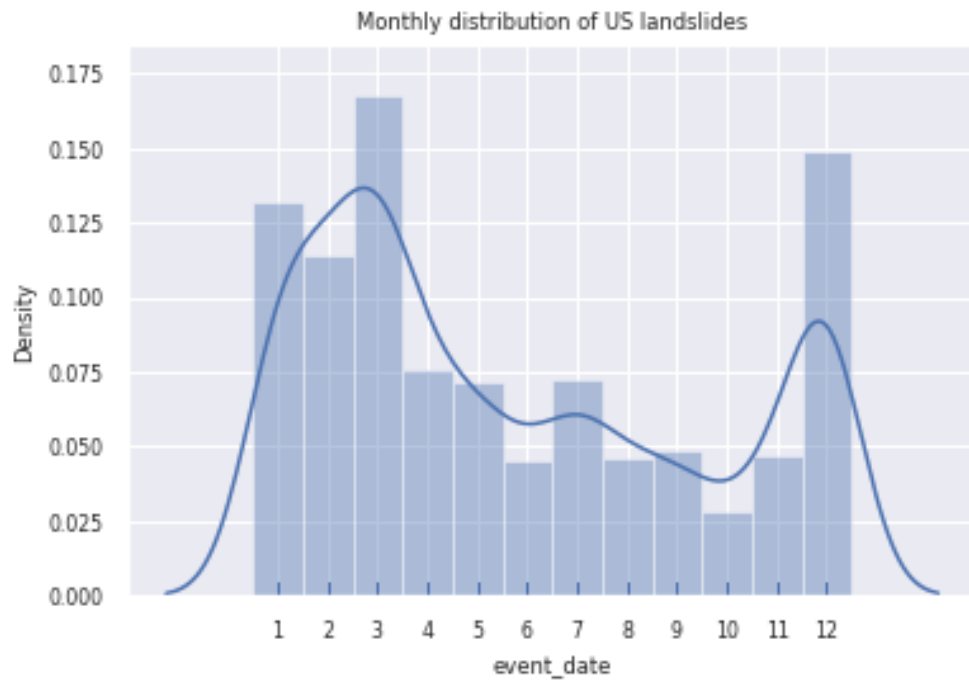


- **Phân bố hàng tháng các lần xảy ra sạt lở đất ở một số quốc gia**

Tháng diễn ra sạt lở nhiều nhất:

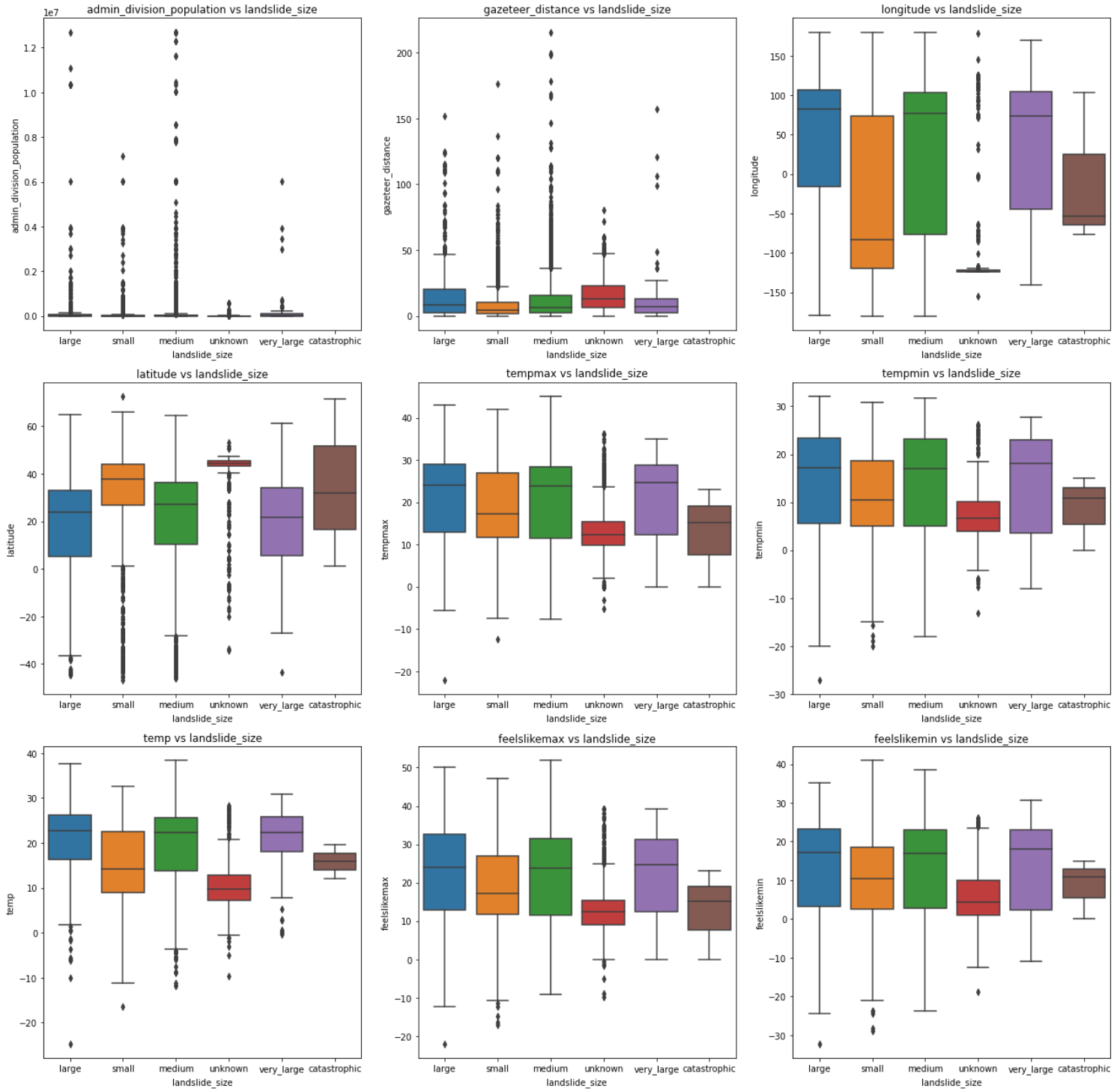
- Tại Mỹ: từ tháng 12 - tháng 3
- Tại Ấn Độ: tháng 6 - tháng 8

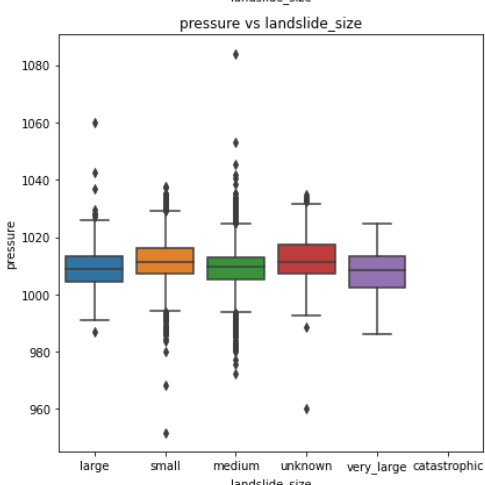
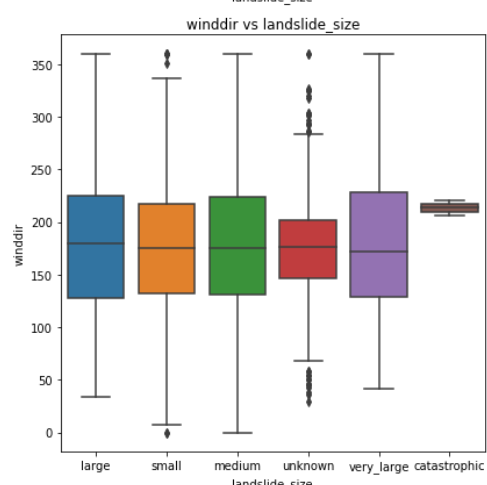
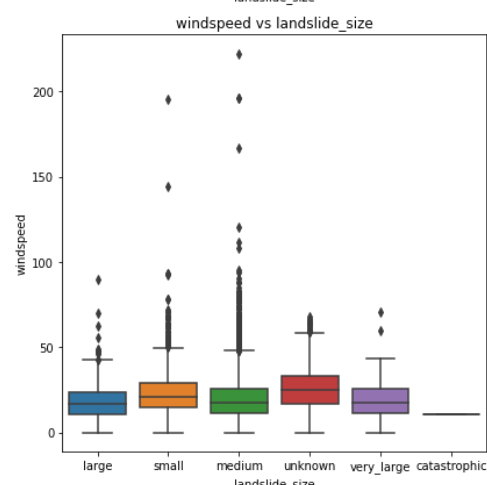
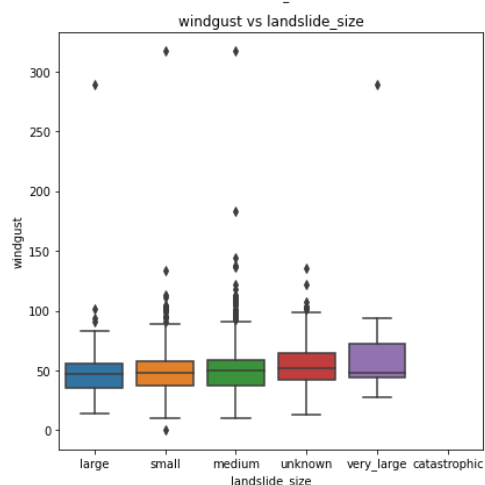
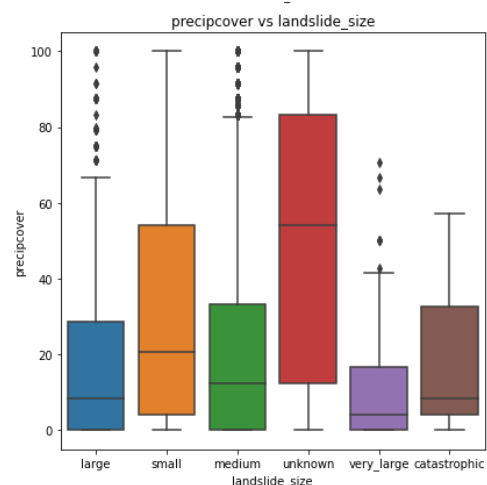
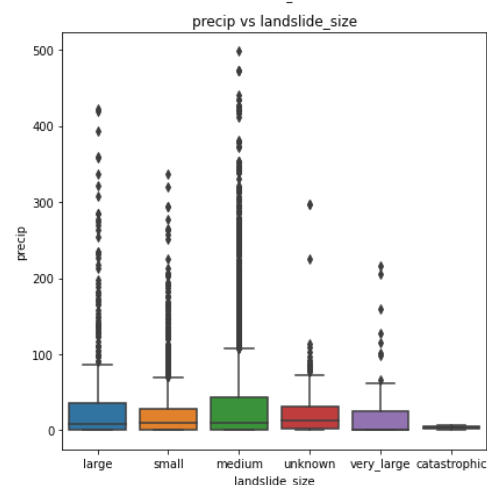
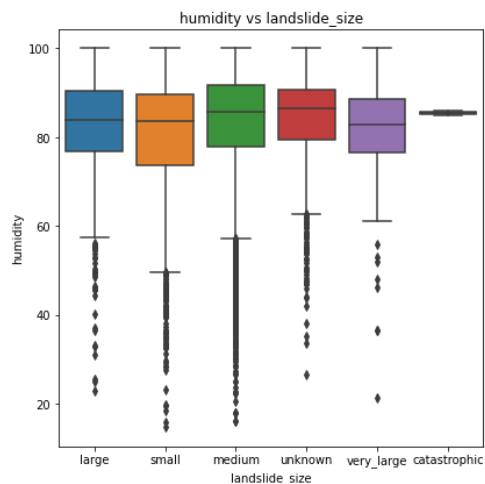
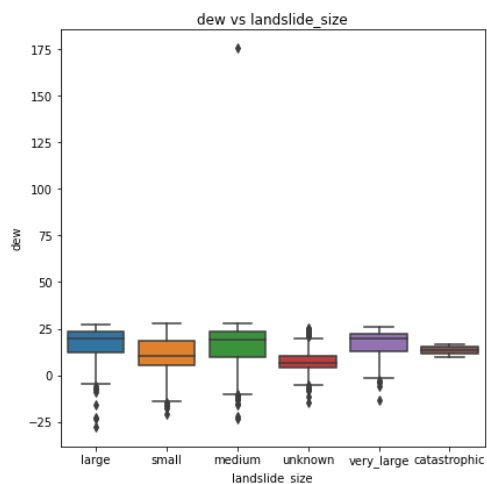
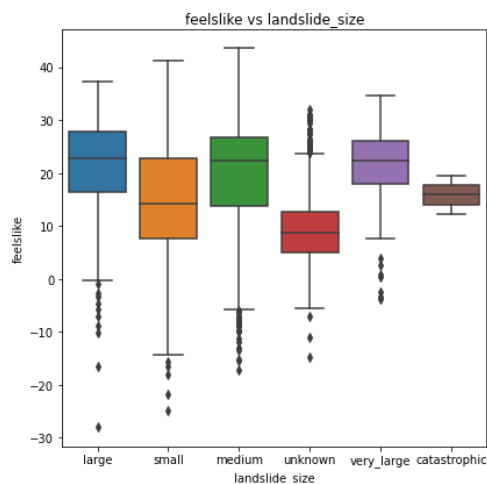
- Tại Việt Nam: tháng 7 - tháng 11
- ⇒ Diễn ra vào mùa mưa bão, bão tuyết



4. EDA dữ liệu thời tiết

- Boxplot dữ liệu thời tiết so với thuộc tính landslide_size





5. Trực quan hóa bản đồ

- Trực quan hóa các địa điểm ghi nhận sạt lở trên thế giới

Phân bố các địa điểm xảy ra sạt lở trên thế giới từ năm 1988 đến năm 2017



- Trực quan hóa các địa điểm ghi nhận sạt lở ở khu vực Châu Á

