

# AN INTEGRATED APPROACH IN IMPLEMENTING OF AN UNSUPERVISED DATA-DRIVEN SYSTEM FOR LANDSLIDES PREDICTION

**Autors:** Luigi Passariello (\*), Maiorana Angela (\*\*\*) . Marco Colombo (\*\*\*) , Stefano Iannello (\*\*\*) , Giuseppe Passariello (\*\*), Fabiano Rinaldi (\*\*\*)

\*INGV: National Institute of Geophysics and Volcanology, public scientific centre of research of Italian Ministry of Research

\*\* Ma.Pa.COM: ICT Company committed to the development of high added value solutions based on 4.0 technologies

\*\*\* CRSLaghi: Lake research and studies centre

**Keywords:** Deep Learning, Landslides prediction, geomorphological data, geological data, climatic data.

## Abstract

We have developed a landslide prediction system, based on the integration of geomorphological, geological and climatic information. The approach to developing the system was to make forecasts using slowly varying factors (geomorphological parameters) and factors with high seasonal variability (soil humidification parameters and rainfall quantities). In this sense, the system involves integration with various real-time data sources such as precipitation forecasting systems, rain gauges and SRS systems. Our objective is to estimate the landslide risk based on the parameters provided as input to the system according to the scale 1. Very low, 2. Low, 3. Medium, 4. High, 5. Very high.

## INTRODUCTION

Historically, both natural and man-made disasters have significant impacts on society, the environment and the economy. The term "disaster" means any situation that has or may have a serious impact on people, the environment or property, including cultural heritage. Over the past 10 years, some high-impact, low-probability (HILP) events such as landslides and floods, caused by forces of nature influenced by climate change, have been identified as the highest risk to society. It should be noted that there is no universally accepted definition of what a low-probability, high-impact risk is. Within the EU, an "operational" definition of these risks has been adopted in order to respond to extreme situations characterized by:

- a. the unpredictability or extraordinary nature of a disaster;
- b. the scale of a disaster, including mass casualties, mass deaths and mass displacement;
- c. the prolonged duration of a disaster;
- d. the degree of complexity of a disaster;
- e. the potential risk of seriously disrupting the functioning of the national government, including the provision of social, environmental, economic and public health services or disruption of critical infrastructure;
- f. geographical area, including the possibility that impacts spread beyond borders;
- g. other factors such as the full activation of the Council's integrated political response mechanisms to crises.

It is also true that, although significant efforts have been made in developing solutions for managing the risk of disasters such as landslides and floods, these have so far proven to be inefficient. Thus, more than 1.5 billion people have been affected by environmental disasters in recent years, with total economic losses exceeding \$1.3 trillion. On the other hand, some of the deadliest disasters in world history have been anthropogenic hazards (environmental disasters caused by human activity) involving nature, these were originally caused or amplified by human instigation resulting in ecological disasters that have changed our world forever.

These risks and associated hazards and vulnerabilities are expected to increase further in the future, due to ongoing changes in land use and demographics, development in hazardous areas, and associated climate change. However, unresolved vulnerabilities of this nature increase both the frequency and magnitude of climate-related disasters and natural hazards resulting in catastrophic loss of life and disruption of livelihoods that inevitably introduce widespread social, environmental and economic impacts. Furthermore, a large number of existing and new risks are interacting in a rapidly changing world, and new correlations between risks are dramatically accelerating climate change and its impacts in a vicious cycle pattern. In this context, the Sendai Framework for Disaster Risk Reduction (DRR) 2015-2030 has identified the need for drastic means to manage the risks of small- and large-scale, frequent and infrequent, sudden and slow-onset disasters caused by natural factors or man-made hazards, as well as related environmental, technological and biological hazards. In particular, in recent years we have seen the emergence of new risks arising from forces of nature or human-triggered events causing mass casualties. Although a large percentage of these hazards occur sporadically or even rarely caused accidentally or intentionally, they have the potential to introduce extreme risks to the population, agriculture, industrial activities, the healthcare system and society. The highly heterogeneous nature of these high-impact, low-probability incidents requires continuous readiness, adequate preparation and planning and the use of advanced technological tools capable of supporting decision-makers in preparation and Early Warning. However, as these hazards are unknown in advance, they require a different mode of operation than routine activities targeting well-known and understood hazards and therefore represent a very complex challenge. For example, what is commonly and improperly defined as a "water bomb" (the correct term would be storm), can cause landslides in places where they had never occurred previously and often also in places where a hydro risk had not been assessed. -high geological. In other cases, despite knowing the risk of flooding or landslides, there are objective difficulties in managing an emergency both in terms of communication with a good part of the population and in terms of evacuating them from the risk areas. Although the threat of landslides and floods, increasingly emerging with ongoing climate change, has a low probability of occurring, if it does occur it can clearly and dramatically affect the economy and security of society, bringing devastation in terms of human life, property, environment and cultural heritage. The effects of climate change on the hydrological cycle, already quite evident today (Intergovernmental Panel on Climate Change, 2022), will presumably continue to manifest themselves in the immediate future with ever greater frequency and intensity. Analyzing the possible scenarios is a necessary operation in order to choose the most effective countermeasures to safeguard the environment and protect living beings. Therefore, a) new combinations of parameters for landslide prediction and b) an integrated approach that combines the use of geomorphological data with climate data will be tested.

## **MATERIALS AND MODELS**

### *Dataset*

With this study we used some data relating to landslide events at an international level and analyzed them using artificial intelligence methods to study the correlation between these events and some geological and climatic parameters. We started from the study of some parameters used to landslide prediction is available in international bibliographies and some evaluations have been made on their total or partial use for an innovative approach that would allow their interpretation related to meteorological events.

Dimension of dataset without the landslide variable is of 2814 rows of 10 columns. Each row presents 10 values (same number of columns of course) with values related to following parameters:

- Topography
- Slope
- Aspect
- Curvature Profile
- Plane of Curvature
- Average rate (occurrence of landslide events)
- Wind
- Humidity Rate
- Solar
- Precipitation (value obtained predictively from meteorological systems or obtained in real time from pluviometric systems and the much more advanced SRS - Smart Rainfall System)

Rainfall, solar exposure, slope and wind all affect the occurrence of landslides. All of these influencing parameters are related to slope aspect [21,22].

Slope. Under normal conditions, rocks and soil, even with sloping slopes, do not spontaneously tend to collapse or slide downhill because the friction with the underlying layers opposes the internal cohesion forces. Landslides occur when the angle of repose is exceeded, which represents the maximum slope beyond which loose materials lose stability and start moving. However, the slope can be a contributing factor, in conjunction with other geomorphological characteristics and the presence of water, to facilitate the triggering of a landslide.

#### Precipitation/Rainfall

Rainfall is an extrinsic variable widely used in susceptibility analysis, and its spatial distribution of annual rainfall is commonly considered in statistical hazard analysis [23, 24]. In study area, landslide was triggered by heavy rainfall. Therefore, annual rainfall data of eight stations around the study area were used to prepare rainfall intensity map.

#### Humidity Rate

Precipitation intensity thresholds are generally calculated assuming that all events are unaffected by antecedent soil moisture conditions. However, it is expected that antecedent soil moisture conditions can provide critical support for the correct definition of trigger conditions. Therefore, we explored the role of antecedent soil moisture on critical thresholds of precipitation intensity and duration to evaluate the possibility of modifying or improving traditional approaches.

#### Slope aspect

It has relation with sunlight exposure, winds, soil moisture content on a slope, and these factors indirectly cause landslide occurrence [25, 26]..

#### Slope curvature (profile and plane)

It is an important geomorphic index of topographic feature, defined as the rate of change of slope in certain direction. It adversely affects the surface erosion by converging and diverging the runoff down the slope [27,28,29].. It has two utmost values: Positive values indicate the surface is convex in upward direction at certain location, while negative values specify that surface is concave in upward direction. Higher the negative value increase, the more the probability of landslide occurrence; on the contrary, comparatively flat area is less exposed to landslide. Slope with concave surface will tend to hold more rainfall water; thus, it has more time for water infiltration into slope and thus increases the probability of landslide occurrence, but the case is opposite for convex slope [28]. The combination of *plan* and *profile* curvature along with curvature is taken into consideration so that slope morphology and flow can be better assessed.

#### *Topographic wetness index*

TWI is used to quantify the topographic control on hydrological process. TWI can measure the degree of accumulation of water at a site. TWI and landslide have a direct relationship if the value of TWI increases the occurrence of landslide [30]. TWI was extracted from DEM.

Average Frequency ratio

FR is a bivariate geo-statistical method to compute the probabilistic correlation between independent and dependent variables [31]. For landslide prediction, it is assumed that the conditioning factors that caused landslide in past may be responsible for initiation landslide in future [32, 26]. The main advantage of FR model is that it is very easy to use and obtain the results that are readily intelligible [33]. FR is based on the correlation of landslides and its conditioning factors. FR is the ratio of the landslides to the total study area; in addition, it is also the ratio of landslide and non-landslide area for a given attribute/class of a parameter. Therefore, calculating FR values, the area ratio with landslide to non-landslide was computed for each class of each factor for the whole study area, and then, area ratio of each class of each factor to the whole study area was calculated. Hence, FR value of each class type was obtained by dividing the ratio of landslide to the ratio of study area [33,32,34].

*Models*

For binary recognition of probability of landslides we use K-means model. It is a partitional clustering algorithm in which each cluster is associated with a centroid and each point is associated with the cluster with the closest centroid. The operation of this algorithm requires that the parameter k, indicating the number of clusters to be assigned, is specified by us. The algorithm works following these steps:

- Select k random points as initial centroids.
- We form the k clusters by assigning all the points to the closest centroid.
- We recalculate the centroids of each cluster.
- If the centroid has changed we repeat the assignment of all points to the closest centroid.
- If the centroid has not changed we have identified the cluster.

Distances can be calculated with any of the measures we have already seen in the previous chapters. Typically the algorithm converges at the first iterations, but stopping criteria can be specified based on the number of iterations, or based on a minimum value of the total error or until the elements of each cluster stabilize (do not pass more from one cluster to another).

The algorithm is sensitive to local minima, as the choice of random starting points could lead to different results in the various executions of the algorithm. A solution may be to run the algorithm several times choosing different starting points until a solution that minimizes the error is obtained, or alternatively other algorithms (such as hierarchical clustering) can be used to determine the starting points. Furthermore, it is strongly affected by the presence of outliers, which tend to "move" the centroid of the cluster, increasing the total error. The solution in this case consists in identifying and eliminating outliers during the data preparation phase.

Tabella 1 – K-Means advantages and disadvantages

<b>Advantages</b>	<b>Disadvantages</b>
<ul style="list-style-type: none"> <li>• Simple to interpret and explain</li> <li>• Flexible and easy to adjust</li> </ul>	<ul style="list-style-type: none"> <li>• It is not as sophisticated as recent clustering algorithms</li> </ul>

- It is efficient and provides good clustering

- Does not guarantee to find the optimal number of clusters
- Requires an assumption about the number of clusters

## RESULTS

To achieve the predictive result we implemented an algorithm whose steps are described and discussed below:

1. **Normalization of the values of the variables;** this facilitates processing and reduces errors within the clustering model

$$x = \frac{x - \min(x)}{\max(x) - \min(x)}$$

2. Dimension of dataset without the landslide variable:(rows: 2814, columns: 10)
3. **Hierarchical clustering;** This is an example that represents the form to generate clusters, the bases will be the same for the K-means

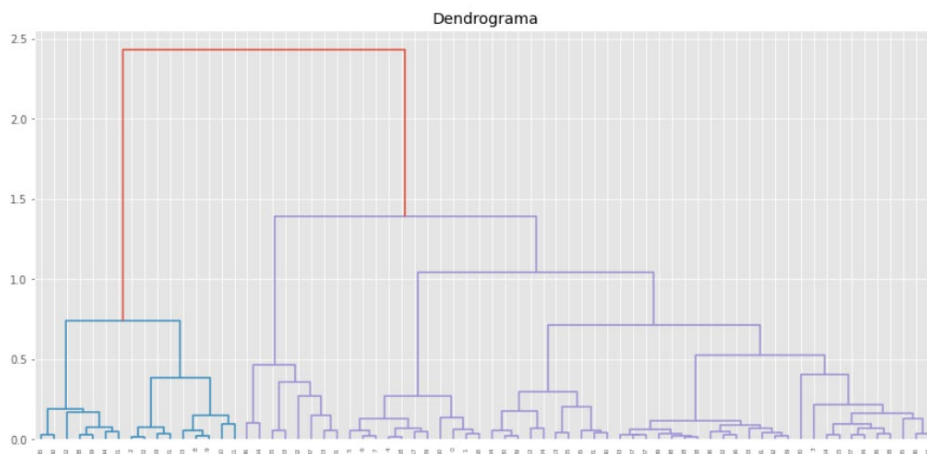


Figure 1 – Dendrogram of Hierarchical clustering

4. **Model Creation;** Determination of ideal number of clusters, for which three graphical methods are used by Elbow method. The elbow method is a technique used in clustering analysis to determine the optimal number of clusters. It involves plotting the within-cluster sum of squares (WCSS) for different cluster numbers and identifying the “elbow” point where WCSS starts to level off. The **elbow method** is a heuristic used in determining the number of clusters in a data set. The method consists of plotting the explained variation as a function of the number of clusters and picking the elbow of the curve as the number of clusters to use. The same method can be used to choose the number of parameters in other data-driven models, such as the number of principal components to describe a data set. In particolare abbiamo utilizzato tre tecniche utilizzate per il metodo ELBOW utilizzando gli algoritmi 1) JAMBÚ, 2) Codo, 3) Silhouette.

- Comparing the results of the three methods we deduced that the main number of the components could be 5 and therefore we applied the K-Kmeans algorithm for k=5
- The summary of the model obtained is as follows:

->Number of Clusters: 5  
 ->Weighting of sample weights: 384.73983341389896  
 ->Numero di etichette: [0 1 2 3 4]  
 ->Numero di iterazioni: 9  
 ->Number of features viewed during the adjustment: 10  
 ->Features seen during the adjustment: ['Topography' 'Slope 2' 'Aspect' 'Curvature Profile' 'Curvature Plane' 'Average' 'Wind' 'Humidity' 'Solar' 'Precipitation']  
 ->Average Coefficient Value: 0.4263854160170544

- The following table outlines for just 20 samples how they are distributed across the various clusters (penultimate row on the right of the table) and which combinations are connected to the occurrence of a landslide or not (last row on the right of the table). In the figure following the table we can see the result of the clustering represented in 3D with only three components. Obviously it is possible to carry out the same visualization considering other combinations of components.

Tabella 2 – How samples are distributed across the various clusters

Sample	Topography	Slope	Aspect	Curvature profile	Plane profile	average Rate	Wind	Humidity	Solar	precipitation	Cluster number	fLandslide
1191	0.012697	0.009595	0.624997	0.542839	0.546710	0.899355	0.001283	0.996546	0.001913	0.997339	2	0
783	0.259759	0.021648	0.719006	0.550882	0.487668	0.967742	0.133168	0.775016	0.186874	0.742398	2	0
1116	0.012411	0.032210	0.695986	0.542842	0.546718	0.898065	0.001018	0.997169	0.001426	0.997968	2	0
1505	0.452227	0.399148	0.118707	0.584182	0.611237	0.477419	0.550213	0.483606	0.554431	0.229467	4	1
2460	0.226612	0.937966	0.479992	0.567740	0.505637	0.865806	0.151738	0.822074	0.175327	0.719015	3	1
2249	0.288143	0.462365	0.966681	0.519721	0.347108	0.601290	0.425099	0.610284	0.430857	0.328729	4	1
520	0.252058	0.069812	0.961786	0.542654	0.544612	0.960000	0.137837	0.774223	0.189900	0.735962	2	0
2578	0.230060	0.784117	0.452453	0.537422	0.477569	0.918710	0.124022	0.822470	0.158841	0.763010	3	1
1956	0.272146	0.097471	0.080438	0.564698	0.624095	0.837419	0.219799	0.743870	0.251287	0.616632	0	1
1771	0.259553	0.534483	0.988247	0.614515	0.971330	0.869677	0.152324	0.818450	0.177379	0.718096	3	1
487	0.267681	0.033253	0.250661	0.543623	0.551925	0.978065	0.124621	0.780169	0.179535	0.754832	0	0
1357	0.740189	0.061098	0.364039	0.542837	0.546718	0.023226	0.979587	0.021745	0.979967	0.009701	1	0
2438	0.898713	0.995042	0.764992	0.537940	0.520476	0.003871	0.996670	0.003907	0.996626	0.001016	1	1
2057	0.261542	0.088459	0.591919	0.561677	0.547022	0.841290	0.217752	0.743530	0.250209	0.619959	2	1
2648	0.192636	0.765177	0.459881	0.535873	0.536864	0.922581	0.123936	0.819752	0.160058	0.762755	3	1
2799	0.218314	0.883962	0.436908	0.597333	0.862724	0.881290	0.144165	0.820998	0.171327	0.731463	3	1
250	0.741918	0.030568	0.386347	0.539888	0.538942	0.023226	0.978846	0.022481	0.979271	0.010100	1	0
901	0.027205	0.007023	0.342758	0.542953	0.564803	0.904516	0.000487	0.993828	0.002782	0.997811	0	0
1347	0.737331	0.132389	0.438561	0.542836	0.546708	0.021935	0.980203	0.021179	0.980558	0.009253	1	0

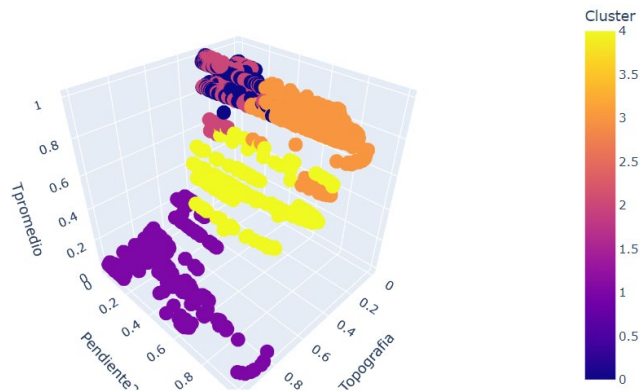


Figure 2 – 3D Representation of Clusters

8. We build a cross-tabulation table that can show the frequency with which certain groups of data appear.

*Tabella 3 – Data Frequency of landslides parameter*

	Cluster	0	1	2	3	4	all
Landslide NO	0	299	601	507	0	0	1407
Landslide YES	1	160	169	210	656	212	1407
		<b>459</b>	<b>770</b>	<b>717</b>	<b>656</b>	<b>212</b>	<b>2814</b>

9. Previous table allows us to make a prediction by association using the created clusters. To carry out landslide prediction tests we can set a combination of values for the parameters considered and use these values to compare them with those obtained in the training phase

Topography = 0.988481

Slope = 0.092862

Aspect = 0.842233

Curvature Profile = 0.547631

Plane of Curvature = 0.582371

Average rate (occurrence of landslide events) = 0.027976

Wind = 0.964126

Humidity Rate = 0.025139

Solar = 0.974434

Precipitation/Rainfall = 0.012548

This combination leads to a "very low" risk of a landslide event. By varying the parameter values, various risk profiles can be obtained.

10. By carrying out 100 manual tests and 200 with automatic choice of parameters within the MinMax range of each parameter, we obtained an accuracy of recognition of conditions that lead to a landslide of 94.7% and recognition of conditions that do not lead to a landslide of 94%. ,9%.

### **PCA – Principal Component Analysis**

Principal component analysis (PCA) is a statistical technique for dimension reduction. In practice, it is used when within a dataset there are many variables correlated with each other and one would like to reduce their number by losing the least amount of information possible.

The principal component analysis (PCA) has the objective of maximizing the variance, calculating the weight to be attributed to each starting variable in order to concentrate them in one or more new variables (called principal components) which will be a linear combination of the starting variables. Maximization of the explained variance, i.e. finding the best combination among all the possible linear combinations of the original variables, i.e. the one that reproduces the greatest portion of variability in the matrix; Orthogonality requirement, i.e. absence of correlation between the components.

Nel nostro caso l'applicazione dell'analisi PCA ci consente di ottenere i seguenti vantaggi:

- **Improve training speeds as the continuously growing data is updated).** The data compressed by PCA provides the important information and is much more digestible by a machine learning model, which now bases its learning on a reduced number of features rather than on all the features present in the original dataset.

- **Feature selection. PCA is essentially a feature selection tool.** When we go to apply it, we look for the features that best explain the variance of the dataset. You can create a ranking of the principal components and order them by importance, with the first component explaining the most variance and the last component explaining the least. By analyzing the principal components it is possible to trace the original features and exclude those that do not contribute to preserving the information in the reduced dimensional plane created by the PCA.
- **Identification of anomalies.** PCA is often used in anomaly identification because it can help identify patterns in data that are not easily distinguishable with the naked eye. Anomalies often appear as data points far from the main group in lower dimensional space, making them easier to detect.
- **Signal identification.** In contrast to the identification of anomalies, PCA is also very useful for signal detection. In fact, just as PCA can highlight anomalies, it can also remove "background noise" that does not contribute to the total variability of the data.

Principal component analysis, or PCA, is a dimensionality reduction method that is often used to reduce the dimensionality of large data sets, by transforming a large set of variables into a smaller one that still contains most of the information in the large set.

Reducing the number of variables of a data set naturally comes at the expense of accuracy, but the trick in dimensionality reduction is to trade a little accuracy for simplicity. Because smaller data sets are easier to explore and visualize and make analyzing data points much easier and faster for machine learning algorithms without extraneous variables to process.

So, to sum up, the idea of PCA is simple: **reduce the number of variables of a data set, while preserving as much information as possible.**

Principal component analysis can be broken down into five steps. I'll go through each step, providing logical explanations of what PCA is doing and simplifying mathematical concepts such as standardization, covariance, eigenvectors and eigenvalues without focusing on how to compute them.

1. Standardize the range of continuous initial variables
2. Compute the covariance matrix to identify correlations
3. Compute the eigenvectors and eigenvalues of the covariance matrix to identify the principal components
4. Create a feature vector to decide which principal components to keep
5. Recast the data along the principal components axes

Il nostro punto di partenza è fornito nella seguente figura:

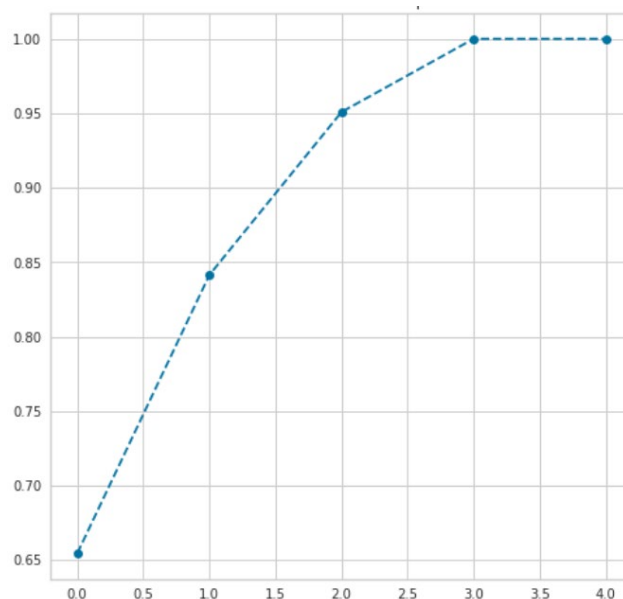




Figure 3 – Variance( axes y) by components (axes x)

**STEP 1: STANDARDIZATION**

The aim of this step is to standardize the range of the continuous initial variables so that each one of them contributes equally to the analysis.

More specifically, the reason why it is critical to perform standardization prior to PCA, is that the latter is quite sensitive regarding the variances of the initial variables. That is, if there are large differences between the ranges of initial variables, those variables with larger ranges will dominate over those with small ranges (for example, a variable that ranges between 0 and 100 will dominate over a variable that ranges between 0 and 1), which will lead to biased results. So, transforming the data to comparable scales can prevent this problem. Mathematically, this can be done by subtracting the mean and dividing by the standard deviation for each value of each variable.

$$z = \frac{\text{value} - \text{mean}}{\text{standard deviation}}$$

Once the standardization is done, all the variables will be transformed to the same scale.

**STEP 2: COVARIANCE MATRIX COMPUTATION**

The aim of this step is to understand how the variables of the input data set are varying from the mean with respect to each other, or in other words, to see if there is any relationship between them. Because sometimes, variables are highly correlated in such a way that they contain redundant information. So, in order to identify these correlations, we compute the covariance matrix.

The covariance matrix is a  $p \times p$  symmetric matrix (where  $p$  is the number of dimensions) that has as entries the covariances associated with all possible pairs of the initial variables. For example, for a 3-dimensional data set with 3 variables  $x$ ,  $y$ , and  $z$ , the covariance matrix is a  $3 \times 3$  data matrix of this form:

$$\begin{bmatrix} Cov(x, x) & Cov(x, y) & Cov(x, z) \\ Cov(y, x) & Cov(y, y) & Cov(y, z) \\ Cov(z, x) & Cov(z, y) & Cov(z, z) \end{bmatrix}$$

Figure 4 – Covariance matrix

**Covariance Matrix for 3-Dimensional Data.**

Since the covariance of a variable with itself is its variance ( $Cov(a,a)=Var(a)$ ), in the main diagonal (Top left to bottom right) we actually have the variances of each initial variable. And since the covariance is commutative ( $Cov(a,b)=Cov(b,a)$ ), the entries of the covariance matrix are symmetric with respect to the main diagonal, which means that the upper and the lower triangular portions are equal.

What do the covariances that we have as entries of the matrix tell us about the correlations between the variables? It's actually the sign of the covariance that matters:

- If positive then: the two variables increase or decrease together (correlated)
- If negative then: one increases when the other decreases (Inversely correlated)

Now that we know that the covariance matrix is not more than a table that summarizes the correlations between all the possible pairs of variables, let's move to the next step.

**STEP 3: COMPUTE THE EIGENVECTORS AND EIGENVALUES OF THE COVARIANCE MATRIX TO IDENTIFY THE PRINCIPAL COMPONENTS**

Eigenvectors and eigenvalues are the linear algebra concepts that we need to compute from the covariance matrix in order to determine the **principal components** of the data. Before getting to the explanation of these concepts, let's first understand what do we mean by principal components.

Principal components are new variables that are constructed as linear combinations or mixtures of the initial variables. These combinations are done in such a way that the new variables (i.e., principal components) are

uncorrelated and most of the information within the initial variables is squeezed or compressed into the first components. So, the idea is 10-dimensional data gives you 10 principal components, but PCA tries to put maximum possible information in the first component, then maximum remaining information in the second and so on, until having something like shown in the scree plot below.

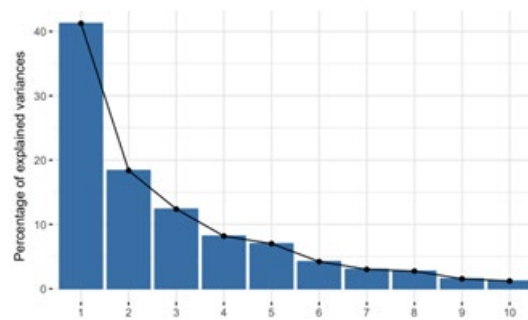


Figure 5 – Principal Components individuation

#### Percentage of Variance (Information) for each by PC.

Organizing information in principal components this way, will allow you to reduce dimensionality without losing much information, and this by discarding the components with low information and considering the remaining components as your new variables.

An important thing to realize here is that the principal components are less interpretable and don't have any real meaning since they are constructed as linear combinations of the initial variables.

Geometrically speaking, principal components represent the directions of the data that explain a **maximal amount of variance**, that is to say, the lines that capture most information of the data. The relationship between variance and information here, is that, the larger the variance carried by a line, the larger the dispersion of the data points along it, and the larger the dispersion along a line, the more information it has. To put all this simply, just think of principal components as new axes that provide the best angle to see and evaluate the data, so that the differences between the observations are better visible.

#### How PCA Constructs the Principal Components

As there are as many principal components as there are variables in the data, principal components are constructed in such a manner that the first principal component accounts for the **largest possible variance** in the data set. For example, let's assume that the scatter plot of our data set is as shown below, can we guess the first principal component? Yes, it's approximately the line that matches the purple marks because it goes through the origin and it's the line in which the projection of the points (red dots) is the most spread out. Or mathematically speaking, it's the line that maximizes the variance (the average of the squared distances from the projected points (red dots) to the origin).

The second principal component is calculated in the same way, with the condition that it is uncorrelated with (i.e., perpendicular to) the first principal component and that it accounts for the next highest variance.

This continues until a total of  $p$  principal components have been calculated, equal to the original number of variables.

Now that we understand what we mean by principal components, let's go back to eigenvectors and eigenvalues. What you first need to know about them is that they always come in pairs, so that every eigenvector has an eigenvalue. And their number is equal to the number of dimensions of the data. For example, for a 3-dimensional data set, there are 3 variables, therefore there are 3 eigenvectors with 3 corresponding eigenvalues. Without further ado, it is eigenvectors and eigenvalues who are behind all the magic explained above, because the eigenvectors of the Covariance matrix are actually *the directions of the axes where there is the most variance* (most information) and that we call Principal Components. And eigenvalues are simply the coefficients attached to eigenvectors, which give the *amount of variance carried*

in each Principal Component. By ranking your eigenvectors in order of their eigenvalues, highest to lowest, you get the principal components in order of significance.

#### STEP 4: FEATURE VECTOR

As we saw in the previous step, computing the eigenvectors and ordering them by their eigenvalues in descending order, allow us to find the principal components in order of significance. In this step, what we do is, to choose whether to keep all these components or discard those of lesser significance (of low eigenvalues), and form with the remaining ones a matrix of vectors that we call *Feature vector*. So, the feature vector is simply a matrix that has as columns the eigenvectors of the components that we decide to keep. This makes it the first step towards dimensionality reduction, because if we choose to keep only  $p$  eigenvectors (components) out of  $n$ , the final data set will have only  $p$  dimensions. So, as we saw in the example, it's up to you to choose whether to keep all the components or discard the ones of lesser significance, depending on what you are looking for. Because if you just want to describe your data in terms of new variables (principal components) that are uncorrelated without seeking to reduce dimensionality, leaving out lesser significant components is not needed.

#### STEP 5: RECAST THE DATA ALONG THE PRINCIPAL COMPONENTS AXES

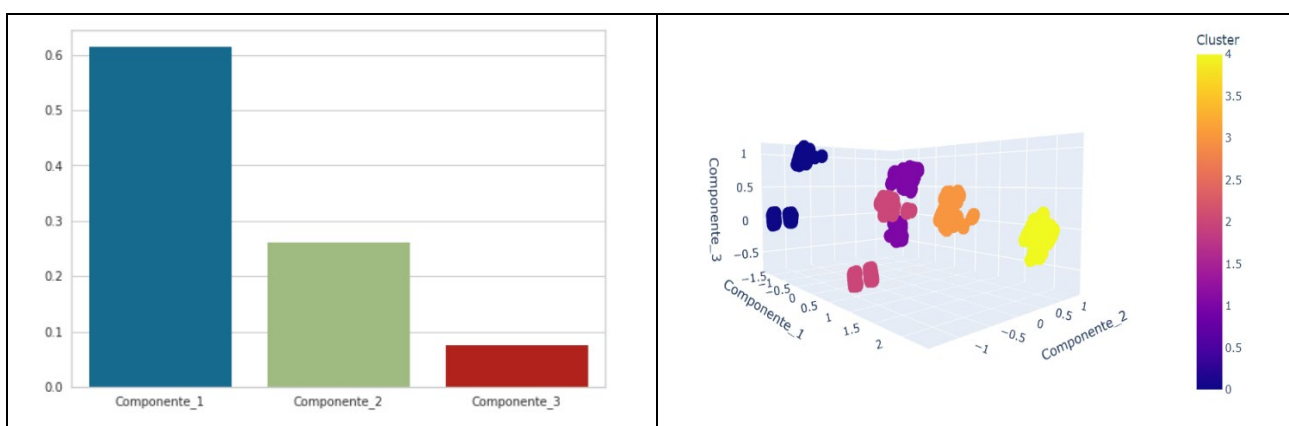
In the previous steps, apart from standardization, you do not make any changes on the data, you just select the principal components and form the feature vector, but the input data set remains always in terms of the original axes (i.e, in terms of the initial variables). In this step, which is the last one, the aim is to use the feature vector formed using the eigenvectors of the covariance matrix, to reorient the data from the original axes to the ones represented by the principal components (hence the name Principal Components Analysis). This can be done by multiplying the transpose of the original data set by the transpose of the feature vector.

$$FinalDataSet = FeatureVector^T * StandardizedOriginalDataSet^T$$

Utilizzando le tre componenti principali abbiamo ottenuto:

Tabella 3 – First 3 main component from PCA Analysis

	Componente_1	Componente_2	Componente_3	Cluster
0	-1.130255	0.957205	-0.330376	1
1	-1.119329	0.963970	-0.310399	1
2	-1.116588	0.957995	-0.445402	1
3	-1.171454	0.966222	-0.230097	1
4	-1.142543	0.965185	-0.301559	1
...	...	...	...	...
2809	1.462817	-0.246777	0.304066	3
2810	1.461691	-0.248802	0.294468	3
2811	1.465988	-0.242886	0.324164	3
2812	1.463248	-0.245803	0.315186	3
2813	1.461042	-0.248161	0.306188	3



*Figure 6 - Total variance bar chart*

*Figure 7 - Three-dimensional representation of the distribution of the clusters*

The distribution of the clusters relating to the three components is given by the following 3D graph.

In the new transformation, using only the aggregate data of the main components and replicating the process analyzed previously, we obtain a schematization of the risk on three factors, low, medium high and a correct identification of the factors that predict a landslide equal to 90.43% and a identification of factors that do not predict a landslide equal to 91.55%. Therefore, the elimination of components with low informative value has allowed us to maintain good accuracy in predicting whether a landslide will occur or not, albeit with the loss of some points of accuracy.

## DISCUSSION

Intervention strategies for landslide prediction. can derive from the use of two distinct approaches: top-down and bottom-up. The top-down approach uses the results of climate projections as input to the different impact models for assessments on the area of interest [36]. However, on the basis of this procedure, the uncertainties inherent in climate models inevitably affect the assessments of the consequences, leading as a precaution to adopt planning strategies that are often too conservative. Furthermore, the continuous evolution of climate projections requires continuous updating of the analyses. On the contrary, the starting point of the bottom-up approach, commonly known as decision scaling, is the assessment of the local vulnerabilities of the system. Climate projections are then used to evaluate their influence on specific local factors and estimate the effectiveness of the chosen strategies. In this regard, the bottom-up approach involves three main phases:

- a. identification of the indicators (and corresponding threshold values) that can best describe the system response, in our case automatically calculable;
- b. construction of a vulnerability domain through a sensitivity analysis of the system to climate variations;
- c. estimate of the changes in the indicators identified as a function of climate projections.

This approach allows for easier identification of suitable solutions, while simultaneously reducing the likelihood of adopting overly precautionary measures. The assessment of the impacts of climate change on geo-hydrological risk undoubtedly represents one of the most serious problems to be addressed. In particular, the estimation of the possible consequences on the stability conditions of landslide areas is of crucial importance in order to guide the choice of the most suitable long-term planning strategies. For the case in question, it was decided to use an integrated analysis tool based on both geomorphological data and data linked to precipitation and therefore to the soil humidity rate. The PCA analysis showed that some components appear predominant compared to all the parameters analyzed to determine landslide phenomena. The analyzes estimate a potential reduction in landslide activity by virtue of the seasonality of the precipitation regime and the substantial increase in temperature that regulates atmospheric evaporative demand. The reduction depends on the time horizon and the socio-economic scenario considered.

The solution, based on the available data, is able to provide a value that identifies the conditions that can trigger a landslide event in terms of differentiated risk conditions: **1. Very low, 2. Low, 3. Medium, 4. High, 5. Very high**. Compared to current models which are mostly based on geological and morphological analyzes which can only explain some of the conditions that predispose to landslides, the study aims to also integrate other conditions such as climatic ones and in particular the humidity on the ground induced by rainfall and their intensity. This parameter allows both an a posteriori evaluation, i.e. after the rains have already occurred, and an a priori evaluation, i.e. dependent on the weather forecasting systems. Therefore, the

integration of the predictive system with a meteorological system allows you to exploit the rainfall data and those of SRS systems (SMART RAINFALL SYSTEM) in real time and predict the probability of a landslide occurring in a given area. In particular, the SRS system (patented by Artys S.r.l.) calculates rainfall maps in real time by analyzing the attenuation of satellite television signals provided by a set of peripheral microwave sensors (35).

## CONCLUSION

The landslide prediction system is based on the integration of geomorphological, geological and climatic information. Since the latter are predictable with forecasting systems, the approach to developing the system was to make forecasts using slowly varying factors (geomorphological parameters) and factors with high seasonal variability (soil humidification parameters and rainfall quantities which affect such factors). In this sense, the system involves integration with various real-time data sources such as precipitation forecasting systems, rain gauges and SRS systems. The latter allow, with a limited number of systems, to define rain intensity maps in real time over large areas. With rain gauges, on the other hand, a very large network of devices would be necessary with the associated maintenance problems that often affect the development but above all the use of a monitoring system in operational contexts. At present the system provides integration with a single source of precipitation data. For a better performance of the system, a broader integration should be envisaged which also involves weather forecasting systems. This may be the subject of future work.

## BIBLIOGRAPHY

- [1] Mondini, A.C., Guzzetti, F. & Melillo, M., *Nat Commun* 14, 2466 (2023). <https://doi.org/10.1038/s41467-023-38135-y>,  
[Google Scholar](#)
- [2] Berg P, Feldmann H, Panitz HJ (2012) Bias correction of high resolution regional climate model data. *J Hydrol* 448–449:80–92. <https://doi.org/10.1016/j.jhydrol.2012.04.026>
- [3] Brunetti MT, Peruccacci S, Antropico L, Bartolini D, Deganutti AM, Gariano SL, Iovine G, Lucani S, Luino F, Melillo M, Palladino MR, Parise M, Rossi M, Turioni L, Vennari C, Vessia G, Viero A, Guzzetti F (2015) Catalogue of rainfall events with shallow landslides and new rainfall thresholds in Italy. In: Lollino G, Giordan D, Crosta GB, Corominas J, Azzam R, Wasowski J, Sciarra N (eds) *Engineering geology for society and territory. Volume 2—landslide processes*. Springer, Berlin, pp 1575–1579
- [4] Coe J, Michael J, Crovelli RA, Laprade WT, Nashem WD (2004) Probabilistic assessment of precipitation-triggered landslides using historical records of landslide occurrence, Seattle, Washington. *Environ Eng Geosci* X:103–122
- [5] De Vita P, Reichenbach P (1998) Rainfall-triggered landslides: a reference list. *Environ Geol* 35:219–233. <https://doi.org/10.1007/s002540050308>
- [6] Gariano SL, Brunetti MT, Iovine G, Melillo M, Peruccacci S, Terranova O, Vennari C, Guzzetti F (2015) Calibration and validation of rainfall thresholds for shallow landslide forecasting in Sicily, Southern Italy. *Geomorphology* 228:653–665. <https://doi.org/10.1016/j.geomorph.2014.10.019>
- [7] Gudmundsson L, Bremnes JB, Haugen JE, Engen-Skaugen T (2012) Technical note: downscaling RCM precipitation to the station scale using statistical transformations—a comparison of methods. *Hydrol Earth Syst Sci* 16:3383–3390. <https://doi.org/10.5194/hess-16-3383-2012>
- [8] Guzzetti F, Peruccacci S, Rossi M, Stark C (2007) Rainfall thresholds for the initiation of landslides in central and southern Europe. *Meteorol Atmos Phys* 98:239–267. <https://doi.org/10.1007/s00703-007-0262-7>

- [9] Reichenbach P, Cardinali M, De Vita P, Guzzetti F (1998) Regional hydrological thresholds for landslides and floods in the Tiber River Basin (central Italy). *Environ Geol* 35:146–159. <https://doi.org/10.1007/s002540050301>
- [10] Segoni S, Piciullo L, Gariano SL (2018) A review of the recent literature on rainfall thresholds for landslide occurrence. *Landslides*. <https://doi.org/10.1007/s10346-018-0966-4>
- [11] Lorenzo Sangelantoni, Eleonara Gioia, Fausto Marincioni, Impact of Climate change on landslides frequency: the Esimno river basin case study (Central Italy) (2018), Springer Science+Business Media B.V, part of Springer Nature, 2018
- [12] Brown C., Meeks R., Hunu K., Yu W. (2010). Hydroclimatic risk to economic growth in Sub-Saharan Africa. *Climatic Change*, DOI: 10.1007/s10584-010-9956-9.
- [13] Comegna L. (2005). Proprietà e comportamento delle colate in argilla. PhD Thesis, Seconda Università di Napoli.
- [14] Cotecchia V., Del Prete M., Federico A., Fenelli G.B., Pellegrino A., Picarelli L. (1984). Some observations on a typical mudslide in a highly tectonized formation in Southern Apennines. *Proc. IV Int. Symp. on Landslides, Toronto, v.2*, 39-44.
- [15] Di Maio C. (1996). Exposure of bentonite to salt solution: osmotic and mechanical effects. *Géotechnique*, 4, 695-707. 5. Guerriero G. (1995). Modellazione sperimentale del comportamento meccanico di terreni in colata. PhD Thesis, Università di Napoli Federico II.
- [16] Iaccarino G., Peduto E., Pellegrino A., Picarelli L. (1995). Principal features of earthflows in part of Southern Apennine. 11th Europ. Conf. On Soil Mechanics and Foundation Engineering, Copenhagen, 4, 354-359.
- [17] Intergovernmental Panel on Climate Change (2022). *Climate Change 2022: Impacts, Adaptation and Vulnerability*. 12th Session of Working Group II and 55th Session of the IPCC, Sixth Assessment Report, <https://www.ipcc.ch/report/ar6/wg2/>.
- [18] Johnson T.E., Weaver C.P. (2009). A framework for assessing climate change impacts on water and watershed systems. *Environmental Management* 43, 118-134.
- [19] Pellegrino A., Picarelli L., Urciuoli G. (2004). Experiences of mudslides in Italy. In L. Picarelli (ed), *Proc. Int. Work on Occurrence and Mechanisms of Flow-Like Landslides in Earthfills and Natural Slopes*, 14-16 Maggio, Sorrento, 191-206,
- [20] Patron, Bologna. 10. Picarelli L., Urciuoli G., Ramondini M., Comegna L. (2005). Main features of mudslides in tectonized highly fissured clay shales. *Landslides*, 2, n. 1, pp. 15-30. 11. Wilby R.L., Dessai S. (2010). Robust adaptation to climate change. *Weather*, 65 (7), pp.180-185.
- [21] Sadr, M.P.; Maghsoudi, A.; Saljoughi, B.S. Landslide susceptibility mapping of Komroud sub-basin using fuzzy logic approach. *Geodyn. Res. Int. Bull.* 2014, 2, XVI–XXVIII. [[Google Scholar](#)]
- [22] Süzen, M.L.; Doyuran, V. Data driven bivariate landslide susceptibility assessment using geographical information systems: A method and application to Asarsuyu catchment, Turkey. *Eng. Geol.* 2004, 71, 303–321. [[Google Scholar](#)] [[CrossRef](#)]
- [23] Dahal RK, Hasegawa S, Nonomura A, Yamanaka M, Dhakal S, Paudyal P (2008) Predictive modelling of rainfall-induced landslide hazard in the Lesser Himalaya of Nepal based on weights-of-evidence. *Geomorphology* 102(3–4):496–510. <https://doi.org/10.1016/j.geomorph.2008.05.041>
- [24] Dahal RK, Hasegawa S, Nonomura A, Yamanaka M, Masuda T, Nishino K (2008) GIS-based weights-of-evidence modelling of rainfall-induced landslides in small catchments for landslide susceptibility mapping. *Environ Geol* 54(2):311–324. <https://doi.org/10.1007/s00254-007-0818-3>
- [25] Yalcin A (2008) GIS-based landslide susceptibility mapping using analytical hierarchy process and bivariate statistics in Ardesen (Turkey): comparisons of results and confirmations. *CATENA* 72(1):1–12. <https://doi.org/10.1016/j.catena.2007.01.003>
- [26] Rahman G, Rahman A, Ullah S, Miandad M, Collins AE (2019) Spatial analysis of landslide susceptibility using failure rate approach in the Hindu Kush region, Pakistan. *J Earth Syst Sci* 128(3):1–16

- [27] Kayastha P, Dhital MR, De Smedt F (2013) Application of the analytical hierarchy process (AHP) for landslide susceptibility mapping: a case study from the Tinau watershed, west Nepal. *Comput Geosci* 52:398–408. <https://doi.org/10.1016/j.cageo.2012.11.003>
- [28] Lee S, Ryu JH, Min K, Won JS (2003) Landslide susceptibility analysis using GIS and artificial neural network. *Earth Surf Process Landf* 28(12):1361–1376. <https://doi.org/10.1002/esp.593>
- [29] Pradhan B (2010) Application of an advanced fuzzy logic model for landslide susceptibility analysis. *Int J Comput Intell Syst.* <https://doi.org/10.1080/18756891.2010.9727707>
- [30] Lee S, Min K (2001) Statistical analysis of landslide susceptibility at Yongin, Korea. *Environ Geol* 40(9):1095–1113. <https://doi.org/10.1007/s002540100310>
- [31] Oh H-J, Kim Y-S, Choi J-K, Park E, Lee S (2011) GIS mapping of regional probabilistic groundwater potential in the area of Pohang City, Korea. *J Hydrol* 399(3):158–172. <https://doi.org/10.1016/j.jhydrol.2010.12.027>
- [32] Lee S, Pradhan B (2007) Landslide hazard mapping at Selangor, Malaysia using frequency ratio and logistic regression models. *Landslides* 4(1):33–41. <https://doi.org/10.1007/s10346-006-0047-y>
- [33] Bourenane H, Guettouche MS, Bouhadad Y, Braham M (2016) Landslide hazard mapping in the Constantine city, Northeast Algeria using frequency ratio, weighting factor, logistic regression, weights of evidence, and analytical hierarchy process methods. *Arab J Geosci* 9(2):154. <https://doi.org/10.1007/s12517-015-2222-8>
- [34] Ozdemir A, Altural T (2013) A comparative study of frequency ratio, weights of evidence and logistic regression methods for landslide susceptibility mapping: Sultan Mountains, SW Turkey. *J Asian Earth Sci* 64:180–197. <https://doi.org/10.1016/j.jseaes.2012.12.014>
- [35] A.Caridi et al. (2018), A field assessment of a novel rain measurement system based on earth-to-satellite microwave links, Conference: WMO/CIMO 2018 Technical Conference on Meteorological and Environmental Instruments and Methods of Observation At: Amsterdam, Netherland
- [36] R.L.Wilby and S. Dessai, (2010), Robust adaptation to climate change, Royal Meteorological Society, , <https://doi.org/10.1002/wea.543>