

Identification of FEP critical paths from a Bayesian network model for the risk assessment of nuclear waste repositories

Edoardo Tosoni

Politecnico di Milano, Italy. Aalto University, Finland. E-mail: edoardo.tosoni@aalto.fi

Francesco Di Maio

Energy Department, Politecnico di Milano, Italy. E-mail: francesco.dimaio@dpolimi.it

Enrico Zio

Energy Department, Politecnico di Milano, Italy. MINES ParisTech/PSL Universit Paris, Centre de Recherche sur les Risques et les Crises (CRC), Sophia Antipolis, France. Eminent Scholar, Department of Nuclear Engineering, Kyung Hee University, South Korea. E-mail: enrico.zio@dpolimi.it

Bayesian networks can be used for the risk assessment of nuclear waste repositories by (i) modeling the causal relations among the set of Features, Events and Processes (FEPs), such as water flows and chemical concentrations, and (ii) calculating the probability that a safety threshold, e.g., on the radionuclide discharge to the environment, is violated. An important outcome of the safety assessment is also the identification of critical paths (i.e., particular combinations of FEP states) leading to such violations. To address this problem, we propose a recursive unsupervised procedure, based on spectral clustering and fuzzy-c-means, for generating mutually exclusive collectively exhaustive clusters of paths covering the possible system states. Then, the probability of each path conditioned on the violation of the safety threshold is evaluated to identify the most critical paths. The procedure is applied to an illustrative deep geological repository.

Keywords: Nuclear waste repositories, Critical paths to failure, Bayesian Networks, Unsupervised spectral clustering, Fuzzy-c-means, Scenario analysis, Risk management.

1. Introduction

Nuclear power plants generate hazardous waste that must be contained to protect humans, which can be achieved through geological burial of nuclear waste. In view of having a nuclear waste repository licensed for construction and operation, it is necessary to perform a safety assessment. Bayesian Networks (BNs) (Pearl and Russell, 2003) are being considered to support this assessment (Tosoni et al., 2019). The nodes of the network represent the Features, Events and Processes (FEPs) that affect the evolution of the repository and its environment. Arcs from one node to another represent the causal dependencies between the FEPs. A sink node represents the safety target, namely the radiological impact relevant to the safety assessment (e.g., the radionuclide discharge to the environment). The risk of the repository can be evaluated as the probability that the safety target violates some safety threshold value.

A given FEP can assume different states, such as low, medium and high, so that the system state is given by the states of all FEPs jointly. The system states which contribute most to risk can

be identified through risk importance measures (Zio, 2011). Here, as a measure to identify such system states, we consider the probability of a system state conditioned on the occurrence of the violation of the safety threshold.

As the number of possible system states to evaluate is very large, it can be meaningful to partition them into fewer mutually exclusive collectively exhaustive subsets. Particularly, we focus on subsets referred to as “*paths to violation* (of the safety threshold)”, that is, to chain of states from one or more independent FEPs all the way through the BN and, hence, to the safety target. However, there are cognitive and computational advantages only if the paths are chosen in an automated, rather than judgmental, way.

Thus, we present a recursive unsupervised procedure based on spectral clustering and a fuzzy-c-means algorithm (Baraldi et al., 2012, 2013) to partition the possible paths into clusters. The procedure starts from the safety target and proceeds backwards in the BN, thus generating sets of mutually exclusive collectively exhaustive paths. The most critical cluster is identified by the largest probability conditioned on violation of the safety threshold.

The practical benefit of the proposed approach is that analysts obtain the most critical paths to the violation of safety threshold without resorting to their own judgment. The knowledge of these paths can inform strategies for lowering the risk level of the repository.

The paper is structured as follows: after recalling the essentials of BNs for risk assessment (Section 2.1) and defining the paths (Section 2.2), we present the novel procedure to identify the most critical ones (Section 2.3). Results with respect to an illustrative deep geological repository are presented in Section 3. Section 4 concludes the work.

2. Methodology

2.1. Bayesian networks for nuclear waste repositories

Nuclear waste repositories can be modeled as BNs whose nodes represent the FEPs and the safety target and whose arcs represent their causal dependences. As an example, Fig. 1 shows the FEPs (white nodes) *Chloride concentration*, *Groundwater flow* and *Canister breach*, and the safety target *Radionuclide discharge* (black) (Tosoni et al., 2019). The identification of all significant FEPs is a crucial task in scenario analysis for nuclear waste management, which should rely on systematic expert elicitation as in the safety assessment of the Yucca Mountain repository (SNL, 2008). Nevertheless, this aspect lies outside the scope of the present paper, which therefore focuses on the illustrative case of Fig. 1.

Formally, let $V = \{i | i = 1, \dots, n_{FEP}, t\}$ be the set of the n_{FEP} FEPs and the safety target $t = n_{FEP} + 1$. The set $A = \{(j, i) | i, j \in V, i \neq j\}$ includes the directed arcs $(j, i) \in A$, indicating that the state of node i depends on j . The nodes $V_-^i = \{j | (j, i) \in A\}$ with an arc to node $i \in V$ are the *parents* of i , which is their *child*. The sets $V^I = \{i | i \in V, V_-^i = \emptyset\}$ and $V^D = \{i | i \in V, V_-^i \neq \emptyset\}$, with $V^D = V \setminus V^I$, contain the independent and the dependent nodes (without and with parents, respectively). For instance, in Fig. 1 *Chloride concentration* is an independent node, whereas *Radionuclide discharge* depends on *Canister breach* and *Groundwater flow*.

Each node $i \in V$ is associated with a random variable X^i with discrete states $s^i \in S^i$, representing intensities such as *low*, *medium*, *high* (Maio et al., 2015). For the safety target $t \in V$, the state s_{vio}^t indicates the violation of a predefined safety threshold (e.g., the regulatory limit). An independent node $i \in V^I$ is in state $s^i \in S^i$ with probability p_{s^i} . For a dependent node $i \in V^D$, $p_{s^i | s_-^i}$ is the conditional probability of state s^i given that its parents $j \in V_-^i$ are in states

$s_-^i \in S_-^i = \times_{j \in V_-^i} S^j$ (where \times denotes the Cartesian product). Referring to Fig. 1, the state probabilities in Tabs. 1, 2 and 3 are taken from the illustrative example of Tosoni et al. (2019).

Table 1. Unconditional state probabilities for *Chloride concentration* (a, left) and *Groundwater flow* (b, right) in Fig. 1.

(a)		(b)	
Chloride concentration		Groundwater flow	
State		State	
Dilute	0.630	Low	0.570
Brine	0.310	Medium	0.250
Saline	0.060	High	0.180

Table 2. Conditional state probabilities for *Canister breach* (a, left) and conditional probabilities that *Radionuclide discharge* violates the safety threshold (b, right) in Fig. 1.

Canister breach				
Chloride concentration	Parents' states		State	
	Groundwater flow	Damage	Severe	Complete
Dilute	Low	0.750	0.200	0.050
Brine	Low	0.420	0.300	0.280
Saline	Low	0.250	0.290	0.460
Dilute	Medium	0.360	0.340	0.300
Brine	Medium	0.250	0.500	0.250
Saline	Medium	0.100	0.120	0.780
Dilute	High	0.220	0.350	0.430
Brine	High	0.020	0.160	0.820
Saline	High	0.050	0.100	0.850

Table 3. Conditional state probabilities for *Canister breach* (a, left) and conditional probabilities that *Radionuclide discharge* violates the safety threshold (b, right) in Fig. 1.

Radionuclide discharge			
Groundwater flow	Parents' states		State
	Canister breach	Violation	
Low	Damage	0.001	
Medium	Damage	0.005	
High	Damage	0.010	
Low	Severe	0.090	
Medium	Severe	0.100	
High	Severe	0.120	
Low	Complete	0.250	
Medium	Complete	0.300	
High	Complete	0.500	

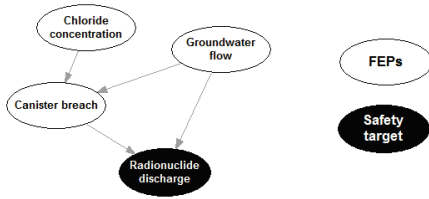


Fig. 1. Bayesian network representing the near field of a deep geological repository.

Now, let \mathbf{s} denote the system state, i.e., a combination of states of all the FEPs so that $\mathbf{s} \in S = \prod_{i \in V^{FEP}} S^i$, where $V^{FEP} = V \setminus t$ is the set of FEPs. The joint probability of each such combination and the violation state of the safety target is

$$p(\mathbf{s}, s_{vio}^t) = \prod_{i \in V^t} p_{s_i} \cdot \prod_{\substack{j \in V^D \\ j \neq t}} p_{s_j | s_{V_j^-}} \cdot p_{s_{vio}^t | s_{V_t^-}}, \quad (1)$$

where s_i and $s_{V_i^-}$ are the states of i and V_i^- as specified by \mathbf{s} , respectively.

Then, the total probability of violating the reference threshold

$$p_{vio} = p(s_{vio}^t) = \sum_{\mathbf{s} \in S} p(\mathbf{s}, s_{vio}^t) \quad (2)$$

is the sum of these joint probabilities over all combinations of FEP states. This *violation probability* can be taken as an estimate of the risk of the repository. In the example of Fig. 1, given the probabilities of Tabs. 1, 2 and 3, the probability that *Radionuclide discharge* violates the safety threshold is 12.5%.

2.2. Critical paths

The overall risk, assessed through the violation probability of Eq. 2, helps decide whether the repository respects the safety requirements. In view of reducing the risk level, however, it is also important to identify which systems states \mathbf{s} contribute to risk most in the sum of Eq. 2.

Rather than evaluating all system states individually, it can be more informative to examine subsets of the set S of system states (preferably mutually exclusive and collectively exhaustive, so that there are no ambiguities and all possible realizations are covered). In what follows, we focus on subsets to which we refer as *paths*.

A FEP can be specified either in terms of occurrence (e.g., *high chloride concentration* or nonoccurrence (e.g., *not low groundwater flow*) of its states. Then, a path specifies the states of a set of FEPs $V^{path} \subseteq V^{FEP}$ such that i) it includes at least one independent node, ii) if it includes a

dependent node, then it also includes at least one of its parents, iii) it includes at least one parent of the safety target.

Essentially, a path captures the chain of states from one or more independent FEPs all the way through the BN and, hence, to the safety target. Here, it can lead to the violation of the safety threshold with some probability. In this sense, paths are analogous to cut sets in but for models in which dependencies are probabilistic (and not described by logical gates) and failure is not defined at the level of system components (e.g., chemical concentrations do not “fail”).

Mathematically, a path can be represented as a vector \mathbf{z} constructed as the concatenation of as many vectors as the FEPs in V^{path} . Each of these vectors, in turn, is composed by as many binary variables as the states of the corresponding FEP $i \in V^{path}$, so that $z^i = 1$ if the state $s^i \in S^i$ is admitted in the path and $z^i = 0$ otherwise. For each FEP, at least one state should be admitted (to rule out impossibility), but not all of them should (as this would be noninformative, and equivalent to not including the FEP in the path), i.e.,

$$1 \leq \sum_{s^i \in S^i} z^i \leq |S^i| - 1, \quad \forall i \in V^{path}. \quad (3)$$

Eventually, each path may lead or not to violating the reference threshold at the safety target. In order to quantify the contribution of a path \mathbf{z} to the probability of such violation, we consider its probability conditioned on violation having occurred (Xu, 2017)

$$\pi(\mathbf{z}) = p(\mathbf{z} | s_{vio}^t) = \frac{p(\mathbf{z}, s_{vio}^t)}{p_{vio}}. \quad (4)$$

The conditional probability of Eq. 4 takes into account both the probability that a path occurs and that it leads to violation once occurred. Therefore, when a set of mutually exclusive collectively exhaustive paths is considered, the sum of each path conditional probability represents the share of violation probability contributed by a path. Accordingly, we identify the most critical path as that with the largest probability conditioned on the occurrence of violation.

2.3. A novel procedure for the identification of critical paths

We propose an automated procedure based on spectral clustering and a fuzzy-c-means algorithm for partitioning the possible system states into informative paths, and identifying the most critical among these. The procedure is also sketched in Fig. 2.

Pre-processing consists of expressing the qualitative labels (such as *low*, *medium*, *high*) for the FEP states $s^i \in S^i, i \in V_t^-$ of the safety target’s

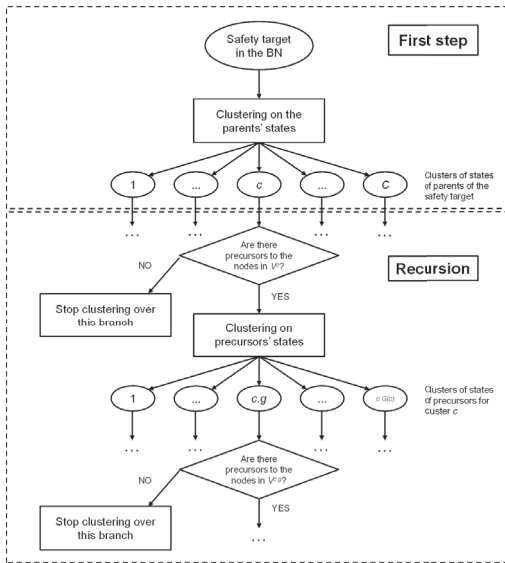


Fig. 2. Automated procedure for generating paths.

parents in quantitative indexes $j = 1, \dots, |S^i|$. Then, data elements $[s^t, p_{s^t_{vio}} | s^t]$ are created for all combinations $s^t \in S^t$ of states of the safety target's parents, so that each data element corresponds to a row of Tab. 3. The rationale is that these data elements describe how different conditions in terms of *Groundwater flow* and *Canister breach* lead to violation with different probabilities.

The unsupervised spectral clustering (details in the Appendix) is, then, launched for partitioning the data elements among the clusters $c = 1, \dots, C$, where C is the number of clusters identified by eigenvalue analysis, as illustrated in the appendix. Each cluster (see Tab. 4 for the data elements related to the parents of the safety target) is identified by specific states of a set of FEPs $V^c \subseteq V^t$ (see Tab. 5).

Hence, for each cluster $c = 1, \dots, C$, we recursively look for clusters of *precursors* states. By precursors, we mean the set of FEPs $V^c = \cup_{i \in V^c} V^i$ of all parents of the FEPs whose states identify the c -th cluster. If this set is empty, there are no precursors to examine. Otherwise, the states $s^j \in S^j, \forall j \in V^c$ should be pre-processed as in the first step. The data elements $[s^c, p_{c|s^c}]$ are created for all combinations s^c of states of the FEPs in V^c . Here, $p_{c|s^c}$ is the conditional probability of the states which identify the c -th cluster, given that the FEPs in V^c are in states s^c . Thus, for cluster $c = 3$ for instance, each data element describes different conditions in terms

Table 4. Clusters characterized by different states of the parents of *Radionuclide discharge* in Fig. 1 and different conditional probabilities of violation.

Radionuclide discharge			
Parents' states		Conditional probability of violation	Cluster
Groundwater flow	Canister breach		
Low	Damage	0.001	3
Medium	Damage	0.005	3
High	Damage	0.010	3
Low	Severe	0.090	3
Medium	Severe	0.100	3
High	Severe	0.120	3
Low	Complete	0.250	2
Medium	Complete	0.300	2
High	Complete	0.500	1

Table 5. FEP states identifying the clusters of states of the safety target's parents.

Cluster c	V^c	States
1	Groundwater flow, Canister breach	High, complete
2	Groundwater flow, Canister breach	Not high, complete
3	Canister breach	Not complete

of *Chloride concentration* and *Groundwater flow* and the corresponding probability of leading to *not complete Canister breach*.

The unsupervised spectral clustering in the Appendix can now be launched to find clusters among these data elements. Each cluster can be identified by specific states of a set of FEPs $V^{c.g} \subseteq V^c, c = 1, \dots, C, g = 1, \dots, G(c)$, where $G(c)$ is the number of clusters for the precursors of cluster c . The obtained clusters describe different modes and chances to result in the c -th cluster.

The example of Fig. 1 illustrates peculiarities that may happen when searching for clusters in BNs. One is that, if a FEP is parent to another FEP in V^c , it appears in both V^c and V^c . For instance, clusters 1 and 2 from the first step are identified by states of *Groundwater flow*, which therefore belongs to V^1 and V^2 . However, being parent of *Canister breach* (the other FEP in V^1 and V^2), it is also in V^1 and V^2 . As it can be seen in Tabs. 6 and 7, this implies the probabilities $p_{c|s^c}, c = 1, 2$, to be null for some combinations of states $s^c, c = 1, 2$, due to inconsistency with the c -th cluster. In these cases, such combinations are excluded from clustering.

Another particular result is that clusters may be identified by multiple sets of FEP states. Consider Fig. 3, which refers to the precursors of cluster 3.

Table 6. Clusters of combinations of states of the precursors of Cluster 1. Inconsistent combinations (i.e., those with *low* or *medium Groundwater flow*) are excluded from clustering.

Cluster 1			
Precursors' states			Cluster
Chloride concentration	Groundwater flow	Conditional probability of Cluster 1	
Dilute	Low	-	-
Brine	Low	-	-
Saline	Low	-	-
Dilute	Medium	-	-
Brine	Medium	-	-
Saline	Medium	-	-
Dilute	High	0.43	1.1
Brine	High	0.82	1.1
Saline	High	0.85	1.1

Table 7. Clusters of combinations of states of the precursors of Cluster 2. Inconsistent combinations (i.e., those with *high Groundwater flow*) are excluded from clustering.

Cluster 2			
Precursors' states			Cluster
Chloride concentration	Groundwater flow	Conditional probability of Cluster 2	
Dilute	Low	0.05	2.2
Brine	Low	0.28	2.2
Saline	Low	0.46	2.1
Dilute	Medium	0.30	2.2
Brine	Medium	0.25	2.2
Saline	Medium	0.78	2.1
Dilute	High	-	-
Brine	High	-	-
Saline	High	-	-

Had the combination *brine, medium* been assigned to cluster 3.1, the identification of clusters 3.1 and 3.3 would imply a unique set of FEPs, as it was the case for the first step in Tab. 5. Instead, there is a sort of nonconvexity of the clusters, which requires that they be identified by multiple sets of FEP states, as shown in Tab. 8.

The recursion proceeds by examining the precursors $V_-^{c.g} = \cup_{i \in V^{c.g}} V_-^i$ of each cluster $c.g, c = 1, \dots, C, g = 1, \dots, G(c)$ from the previous step. Again, the unsupervised spectral clus-

	Low	Medium	High
Dilute	3.3	3.2	3.2
Brine	3.3	3.3	3.1
Saline	3.3	3.1	3.1

Fig. 3. Partition of the states of the precursors of cluster 3 *Chloride concentration* (rows) and *Groundwater flow* (columns). The assignment of the combination *brine, medium* to cluster 3.3 implies a double identification for this cluster and cluster 3.1 (see Tab. 8).

Table 8. FEP states identifying the clusters of states of the precursors of cluster 3.

Cluster 3.g	$V^{3.g}$	States
3.1	Chloride concentration, Groundwater flow	Not dilute, high
	or	
	Chloride concentration, Groundwater flow	Saline, medium
3.2	Chloride concentration, Groundwater flow	Dilute, not low
3.3	Groundwater flow	Low
	or	
	Chloride concentration, Groundwater flow	Brine, medium

tering in the Appendix partitions the data elements $[s_-^{c.g}, p_{c.g|s_-^{c.g}}]$ among a number $H(c.g)$ of clusters $c.g.h, h = 1, \dots, H(c.g)$, which are identified by specific states of a set of FEPs $V^{c.g.h} \subseteq V_-^{c.g}$. The obtained clusters describe different modes and chances to result in the $c.g$ -th cluster. The procedure terminates at a step such that $V_-^{c.g.h\dots k} = \emptyset, \forall c = 1, \dots, C, \forall g = 1, \dots, G(c), \forall h = 1, \dots, H(c.g), \dots, \forall k = 1, \dots, K(c.g.h\dots)$ (where $K(c.g.h\dots)$ is the number of clusters of precursors of cluster $c.g.h\dots$). The number of superscripts reflects the overall number of steps. In our example, we stopped at the second step, as neither *Chloride concentration* nor *Groundwater flow* have parents.

3. Results

Once the recursive procedure has terminated, we can display all identified clusters, starting from those of the safety target parents and proceeding with the various branchings of precursors (Tab. 9). Due to the recursive nature of the novel procedure, the conditions progressively specified

along each branch meet the criteria for constituting paths. Furthermore, because the unsupervised spectral clustering in the Appendix partitions the FEP states at each step, these paths are mutually exclusive and collectively exhaustive.

These paths are ranked in Tab. 9 by their probability conditioned on violation. The most critical one is *high groundwater flow, complete canister breach*, which contributes 41% of the violation probability. It is followed at some distance by *not saline chloride concentration, not high groundwater flow, complete canister breach* (30%), and *low groundwater flow* (or *brine chloride concentration, medium groundwater flow*), *not complete canister breach* (13%).

Table 9. Paths for the repository of Fig. 1 ranked by probability conditioned on the violation of the safety threshold on Radionuclide discharge.

Chloride concentration	Scenario Groundwater flow	Canister breach	Conditional probability
-	High	Complete	0.41
Not saline	Not high	Complete	0.30
-	Low	Not complete	0.13
or	Medium		
Brine	Not low	Not complete	0.09
Dilute	Not high	Complete	0.06
Saline	High	Not complete	0.01
Not dilute			
or	Not low		
Saline			

Arguably, the proposed automated procedure generates more informative paths than, say, subjective approaches. In the latter case, in fact, the detail of the partitions may not be entirely adequate to highlight the FEP states that should be addressed in risk management.

For instance, compared to looking at paths made of one FEP only (Tab. 10), we have broadened the scope of the analysis by emphasizing that also *Canister breach* plays an important role in increasing the violation probability. Indeed, a *high groundwater flow* has a quite low conditional probability (1% or 9% in Tab. 9) if not paired with a *complete canister breach*.

At the opposite end, Tab. 11 examines all system states *s* individually (each such combination is colored as its consistent path of Tab. 9). Here, based on the top-ranked combination, one may unnecessarily focus on the occurrence in which *Chloride concentration* is *dilute*. Even worse, he or she may inappropriately disregard the most critical path *high groundwater flow, complete canister*

Table 10. States of *Groundwater flow* ranked by probability conditioned on the violation of the safety threshold on *Radionuclide discharge*.

Groundwater flow	Conditional probability
Low	0.27
Medium	0.27
High	0.46

breach, as this only ranks tenth when matched with *saline chloride concentration*.

Table 11. System states ranked by probability conditioned on the violation of the safety threshold on *Radionuclide discharge*.

Chloride concentration	Scenario Groundwater flow	Canister breach	Conditional probability
Dilute	High	Complete	1.96e-01
Brine	High	Complete	1.84e-01
Dilute	Medium	Complete	1.14e-01
Brine	Low	Complete	9.92e-02
Dilute	Low	Severe	5.19e-02
Brine	Medium	Complete	4.66e-02
Dilute	Medium	Severe	4.30e-02
Brine	Low	Severe	3.83e-02
Dilute	High	Severe	3.82e-02
Saline	High	Complete	3.68e-02
Dilute	Low	Complete	3.60e-02
Saline	Low	Complete	3.16e-02
Brine	Medium	Severe	3.11e-02
Saline	Medium	Complete	2.82e-02
Brine	High	Severe	8.60e-03
Saline	Low	Severe	7.16e-03
Dilute	Medium	Damage	2.27e-03
Dilute	Low	Damage	2.16e-03
Dilute	High	Damage	2.00e-03
Saline	Medium	Severe	1.44e-03
Saline	High	Severe	1.04e-03
Brine	Medium	Damage	7.77e-04
Brine	Low	Damage	5.95e-04
Brine	High	Damage	8.95e-05
Saline	Low	Damage	6.86e-05
Saline	Medium	Damage	6.02e-05
Saline	High	Damage	4.33e-05

The effectiveness of the procedure in identifying informative paths can be sensitive to the form in which the input (namely, the FEP states and probabilities which constitute the data elements) is fed to it. In particular, while probabilities do not necessarily require any modification, the qualitative FEP states forcibly need to be turned into numerical values. In the pre-processing phase of Section 2.3, such values have been taken to be integer indexes. Whatever the choice, however, it is important to use values that are diverse enough not to lump the data elements into very few clusters (similarly to Tab. 10), but also not so unique that the number of clusters nears that of the data elements themselves (similarly to Tab. 11). In the former case, especially, the risk would be to obtain overly large paths whose risk score is very high just because they cover most of the probability space rather than for having any significant impact.

In summary, safety analyst should focus on the path in which *high groundwater flow* is followed by *complete canister breach*. When optimizing the repository design during the licensing process, risk can be lowered by increasing the resistance of the canister especially in conditions of intense water infiltrations.

4. Conclusions

This paper has tackled the problem of identifying the most critical paths, in terms of FEP states, from a Bayesian network model for the risk assessment of nuclear waste repositories. Risk has been evaluated as the probability of violating a safety threshold on the radiological impact. The most critical paths are those combinations of FEP states with the largest probability conditioned on the occurrence of the safety threshold violation.

Since the paths among which to search for the most critical ones are many, we have proposed an automated procedure based on spectral clustering and a fuzzy-c-means algorithm. This procedure generates a set of mutually exclusive collectively exhaustive paths with a level of detail that can be informative for risk management.

Appendix. Unsupervised spectral clustering

Clustering procedures help partition the elements of a dataset into homogeneous groups named *clusters*. When the number and the characteristics of the clusters are unknown ahead of time, we say that clustering is *unsupervised*. The partition can be performed through a procedure consisting of i) measuring the similarity among the elements in the dataset, ii) spectral analysis and iii) a fuzzy-c-means (FCM) algorithm (Baraldi et al., 2012, 2013).

First, the similarity between the elements $\mathbf{x}_i, i = 1, \dots, N$, in the dataset is measured. Each

element is an instance of a vector of K types of data (e.g., states of different system components). To facilitate the recognition of similarities, these raw data may be pre-processed and replaced by some function (e.g., logarithm) or transform (e.g., wavelet). Let us refer to the pre-processed data as $\mathbf{y}_i, i = 1, \dots, N$. These are collected in the matrix Y , whose generic element is denoted by $y_{ik}, i = 1, \dots, N, k = 1, \dots, K$. The distance between any pair \mathbf{y}_i and $\mathbf{y}_j, i \neq j$, in the K -dimensional space is computed as

$$\delta_{ij} = \sum_{k=1}^K (y_{ik} - y_{jk})^2. \quad (5)$$

Subsequently, these distances are turned into similarities μ_{ij} through suitable functions such as

$$\mu_{ij} = e^{\left(\frac{-\ln(\alpha)}{\beta^2} \delta_{ij}\right)}, \quad (6)$$

so that similarities range from 0 (full dissimilarity) to 1 (full similarity). The larger the ratio $-\ln(\alpha)/\beta^2$ the more stringent the evaluation of similarity. Less smooth functional shapes such as triangular and trapezoidal may be also chosen (Dubois et al., 1988), but the sensitivity of the similarity to the implied discontinuities would need to be assessed. The μ_{ij} constitute the elements of the *similarity matrix* W , symmetric of size N , so that the generic cell indicates the similarity between the element on the row and that on the column.

Second, spectral analysis permits to identify the number of clusters among which to partition the data elements. Towards this end, the *normalized graph Laplacian* matrix can be computed as

$$L_{\text{sym}} = I - D^{-\frac{1}{2}} \cdot W \cdot D^{-\frac{1}{2}}, \quad (7)$$

where I is the identity matrix and D is a diagonal matrix whose diagonal elements result from

$$d_i = \sum_{j=1}^N \mu_{ij}, \quad i = 1, \dots, N. \quad (8)$$

The spectrum of L_{sym} carries the information about a graph whose nodes and arcs correspond to the data elements and their similarities, respectively (Alpert et al., 1999; von Luxburg, 2007; Zhao and Liu, 2007; Zio and Maio, 2010). Thus, upon sorting the eigenvalues $\lambda_i, i = 1, \dots, N$, of L_{sym} in ascending order, the number C of clusters can be taken such that the difference between the eigenvalues λ_C and λ_{C+1} is large relative to those between other consecutive pairs of eigenvalues (von Luxburg, 2007). The eigenvectors $\mathbf{u}_c, c = 1, \dots, C$, corresponding to the first C eigenvalues can, then, be taken as the columns of the *reduced matrix* U of size $N \cdot C$. This is a revisited version of the similarity matrix W such that each element of the dataset is expressed as a vector made of the

i -th component, $i = 1, \dots, N$, of the C eigenvectors. The ability to divide data elements among the clusters can be further enhanced by replacing U with its normalized version T , whose generic element is (von Luxburg, 2007)

$$t_{ic} = \frac{u_{ic}}{\sqrt{\sum_{c=1}^C u_{ic}^2}}, \quad i = 1, \dots, N, \quad c = 1, \dots, C. \quad (9)$$

Finally, the N data elements are partitioned into C clusters through the fuzzy-c-means algorithm (Bezdek, 1981). This algorithm estimates the degree of membership (shortly, membership) m_{ic} of each element $i = 1, \dots, N$ to each cluster $c = 1, \dots, C$. These memberships are calculated via minimization of the objective function

$$J = \sum_{i=1}^N \sum_{c=1}^C m_{ic}^\phi \cdot b_{ic}^2, \quad (10)$$

where the larger $\phi > 1$ the fuzzier the clusters. If some level of ‘‘purity’’ of the clusters is needed (whereby a least membership threshold would be set before assigning a data element to a cluster), then this parameter should be kept not too large.

Moreover,

$$b_{ic}^2 = \sum_{c=1}^C (t_{ic} - \gamma_c)^2 \quad (11)$$

is a measure of the distance between the i -th element and the centre γ_c of the c -th cluster. These centres are vectors of length C , each one resulting from the weighted sum of the data elements

$$\gamma_c = \frac{\sum_{i=1}^N m_{ic}^\phi \cdot \mathbf{t}_i}{\sum_{i=1}^N m_{ic}^\phi}, \quad c = 1, \dots, C. \quad (12)$$

As it can be seen, the weights are related to the memberships (elevated to ϕ), which are calculated as

$$m_{ic}^\phi = \frac{b_{ic}^{-2(\phi-1)}}{\sum_{c=1}^C b_{ic}^{-2(\phi-1)}}. \quad (13)$$

After initializing these memberships to feasible values, the objective function J is minimized by iteratively updating the centres, the square distances and again the memberships. The minimization ensures that the data elements have larger memberships to the clusters they are closer to the centres of. Accordingly, each element can be finally assigned to the cluster to which it has the largest membership. This guarantees that clusters are mutually exclusive. Possibly, assignments may be done only above a minimum membership, in which case the cluster of unassigned data elements should be defined so that the clusters are also collective exhaustive.

References

- Alpert, C., A. Khang, and S. Yao (1999). Spectral partitioning: the more eigenvectors, the better. *Discrete Applied Math* 90, 3–26.
- Baraldi, P., F. D. Maio, and E. Zio (2012). Unsupervised clustering for fault diagnosis. *Proceedings of the Prognostics and System Health Management Conference, Beijing, People’s Republic of China*.
- Baraldi, P., F. D. Maio, and E. Zio (2013). Unsupervised clustering for fault diagnosis in nuclear power plant components. *International Journal of Computational Intelligence Systems* 6(4), 764–777.
- Bezdek, J. C. (1981). *Pattern recognition with fuzzy objective function algorithms*. New York, New York: Springer.
- Dubois, D., H. Prade, and C. Testemale (1988). Weighted fuzzy pattern matching. *Fuzzy Sets and Systems* 28, 313–331.
- Maio, F. D., S. Baronchelli, and E. Zio (2015). A computational framework for prime implicants identification in noncoherent dynamic systems. *Risk Analysis* 35(1), 142–156.
- Pearl, J. and S. Russel (2003). *Bayesian networks, in: M.A. Arbib (Eds.), Handbook of Brain Theory and Neural Networks 2nd edition*. Cambridge, Massachusetts: MIT Press.
- SNL (2008). *Features, Events, and Processes for the Total System Performance Assessment: Analysis. ANL-WIS-MD-000027 REV 01*. Las Vegas, Nevada: SNL.
- Tosoni, E., A. Salo, J. Govaerts, and E. Zio (2019). Comprehensiveness of scenarios in the safety assessment of nuclear waste repositories. *Reliability Engineering and System Safety* 188, 561–573.
- von Luxburg, U. (2007). A tutorial on spectral clustering. *Statistics and Computing* 17(4), 396–416.
- Xu, B. G. (2017). Intelligent fault inference for rotating flexible rotors using bayesian belief network. *Expert Systems with Applications* 36, 816–822.
- Zhao, Z. and H. Liu (2007). Spectral feature selection for supervised and unsupervised learning. *Proceedings of the 24th international conference on machine learning, Corvallis, OR*.
- Zio, E. (2011). *Risk importance measures, in: H. Pham (Eds.), Safety and risk modeling and its applications*. London, United Kingdom: Springer Verlaag.
- Zio, E. and F. D. Maio (2010). A data-driven fuzzy approach for predicting the remaining useful life in dynamic failure scenarios of a nuclear system. *Reliability Engineering and System Safety* 95(1), 49–57.