

# Resiliency Assessment in Distribution Networks Using GIS Based Predictive Risk Analytics

Jonatas Boas Leite, *Member, IEEE*, José Roberto Sanches Mantovani, *Member, IEEE*, Tatjana Dokic, *Student Member, IEEE*, Qin Yan, *Student Member, IEEE*, Po-Chen Chen, *Student Member, IEEE*, and Mladen Kezunovic, *Life Fellow, IEEE*

**Abstract**— A new predictive risk-based framework is proposed to increase power distribution network resiliency by improving operator understanding of the energy interruption impacts. This paper expresses the risk assessment as the correlation between likelihood and impact. The likelihood is derived from the combination of Naive Bayes learning and Jenks natural breaks classifier. The analytics included in a GIS platform fuse together a massive amount of data from outage recordings and weather historical databases in just one semantic parameter known as failure probability. The financial impact is determined by a time series-based formulation that supports spatiotemporal data from fault management events and customer interruption cost. Results offer prediction of hourly risk levels and monthly accumulated risk for each feeder section of a distribution network allowing for timely risk mitigation.

**Index Terms**—Power distribution system, risk assessment, Naive Bayes learning, failure probability, time series, interruption cost, geographic information system (GIS).

## I. INTRODUCTION

THE proposed predictive risk management framework leads to pro-active risk management and effective ranking of risk reduction measures [1]. The weather-based risk assessment provides the spatiotemporal correlation between weather data and historical management data of the power distribution system. Historically, the risk assessment was mainly studied in power transmission system, [2] and [3]. The most recent literature on power distribution system has also focused on risk studies as a central theme [4]–[10].

In [4], historic reliability data reflecting the variation of service continuity indices is utilized to develop probability distribution functions used to illustrate the potential financial risk associated with assigned reward/penalty structure integrated in a performance-based regulation plan for distribution utilities. The histograms of indices, such as system

average interruption frequency index (SAIFI) and duration index (SAIDI), overlap a predefined function that reproduces the reward/penalty regulation policy, predicting the future risks. Instead of evaluating the financial risk, [5] introduces a risk assessment approach that ensures the human safety in power distribution network by determining the intensity of fault current levels that are dangerous for people when stepping on downed conductor and touching poles in a faulted network. The risk analysis employs the Monte Carlo simulation using assumptions of probability distribution functions in the soil resistivity, human body resistance and heart current. Another study presented in [6] analyzes the risk from vaults in the underground power distribution system that can provoke human injuries, monetary compensation, energy unavailability and traffic disruption on streets.

More recent issues involving the penetration of renewable energy resources into power distribution system are also being investigated through the risk analysis approach [7]–[10]. In [7], the correlation between day-ahead and real-time markets is integrated in a reliability and price risk assessment using an energy and pre-dispatch model. Going beyond the short-term market operation, work in [8] investigates the risk-based security of concentrated solar power for mid- and long-term planning horizons. The impact indices are aimed at minimizing steady-state voltage profile variation, assessing the line overload security, and verifying the static and dynamic voltage stability. The four severity continuous functions determine the risk using chronological simulation technique with clustered solar generation patterns for each yearly season. Similarly, [9] assesses the impact of increasing the wind power injection into medium-voltage networks. Investment alternatives taking into account photovoltaic generation, electric vehicles and other new technologies at low-voltage network have been assessed by using the planning framework which determines the risks based on availability, losses and power quality [10].

We have proposed several innovative solutions: a) integration of outage records, historical weather information and fault management events in a risk-based GIS driven proactive management tool; b) implementation of a risk model based on Naive Bayes learning, and classifying the calculated likelihood using Jenks natural breaks where the financial impacts are modeled using the time series-based spatiotemporal formulation, and c) operator visualization of risk prediction and mitigation using GIS interface.

This work was fully supported by the São Paulo Research Foundation – FAPESP (grant: 2015/17757-2), CAPES and CNPq (grant: 305371/2012-6) allowing Dr. Boas Leite to spend a year with Dr. Kezunovic's team at Texas A&M University.

J. B. Leite and J. R. S. Mantovani are with the Electrical Engineering Department, UNESP/FEIS, Ilha Solteira, São Paulo, BRAZIL (e-mail: jonatasboasleite@gmail.com, mant@dee.feis.unesp.br).

T. Dokic, Q. Yan, P.-C. Chen and M. Kezunovic are with the Department of Electrical and Computer Engineering at Texas A&M University, College Station, TX, USA (e-mail: tatjana.djokic@email.tamu.edu, judyqinyan2010@gmail.com, pchen01@tamu.edu, kezunov@ece.tamu.edu).

This paper is organized as follows. Section II presents the risk assessment background by introducing a risk metric in a form of risk matrix. Section III proposes the Jenks natural breaks classification method for defining risk matrix row/column classes. Subsequently, the calculation of failure probability and interruption cost as well as the procedure that obtain the risk matrix are demonstrated. In Section IV, explained concepts involving the proposed risk assessment framework are utilized in the evaluation of a real world distribution network. The conclusions are given in Section V before the references at the end.

## II. RISK ANALYSIS BACKGROUND

The proposed predictive risk analysis can offer anticipation of problems that may, or may not have happened before, in order to assist pro-active risk management strategies. The risk analysis framework aims to minimize the energy interruption impacts and to reduce human intervention by performing two steps sequentially in pre-defined time span. The first one is the risk assessment whereas the second one is the risk mitigation through the real-time control technique. In this framework, the risk assessment, the main of this work, is an important step in quantifying each part of the analyzed problem by calculating the likelihood and impact estimates. Equation (1) presents the quantified risk expressed as expected value of loss, i.e. consequences along time are given by the correlation between the likelihood of event occurrence along time and consequent impacts of each event [11].

$$RISK\left(\frac{conseq}{time}\right) = LIKELIHOOD\left(\frac{event}{time}\right) \times IMPACT\left(\frac{conseq}{event}\right). \quad (1)$$

This correlation is typically obtained by a risk matrix where the risk is ranked in levels, thus matrix elements should be grouped in three levels: the high (H) level is considered unacceptable risk; the medium (M) level is dealt as either undesirable or as acceptable with review; and the low (L) level is treated as acceptable without review. The number of rows and columns of the risk matrix is defined by likelihood and impact categories as demonstrated in Table I.

## III. CLASSIFICATION METHODOLOGY

Since levels and categories represent ranges of continuous values, a clustering methodology is needed to classify the estimated likelihood, impact and risk. The Jenks natural breaks algorithm is a common method in GIS applications able to divide a dataset into a predefined number of homogeneous classes and was originally introduced as a

TABLE I. ROWS AND COLUMNS CATEGORIES OF THE RISK MATRIX.

Rows			Columns		
Categories		Description	Categories		Description
LIKELIHOOD	I	Extremely Unlikely	IMPACT	A	Insignificant
	II	Highly Unlikely		B	Minor
	III	Doubtful		C	Significant
	IV	Very Unlikely		D	Serious
	V	Unlikely		E	Major
	VI	Likely		F	Catastrophic

method for "optimal data classification" because it minimizes the variances within classes by maximizing the variance between classes [12]. One-dimensional values, which are not uniformly distributed, fits perfectly into Natural breaks classification [13], consequently, the well-known k-means clustering is its generalization for multivariate data [14].

The Algorithm 1 describes methodically all steps involved in the procedure for obtaining the class boundaries from an input dataset  $U$  using the Jenks optimization algorithm. At the beginning, the class boundaries are defined by intervals with the same size. Then, the algorithm adjusts the boundaries systematically until the minimization of the sum of the squared deviation from the classes i. e. until the maximization of  $GVF$ , that varies into interval  $[0, 1]$ , is achieved. In this way, the algorithm achieve the class boundaries that produce the maximal similarity to data points in a class.

The Jenks natural breaks optimization performs the central role in the determination of boundaries for each class, i.e. inferior and superior limits for each likelihood and impact category as well as for each risk level. In Fig. 1, the input dataset  $U$  into Jenks optimizer comes from calculations of failure probability and interruption cost. The product of probabilities and costs becomes one-dimensional risk dataset permitting to use again the Jenks optimizer on risk level

### Algorithm 1 Jenks Natural Breaks algorithm.

- 1: Select the input dataset  $U$  to be classified and specify the number of classes,  $NC$ .
- 2: Define the classes' boundaries:  $[INF_j, SUP_j]$  to  $j = 1, 2, \dots, NC$ , where every interval has the same size.
- 3: Calculate the sum of squared deviation of the dataset,  $SD_U$ , using (2):

$$SD_U = \sum (u_i - \bar{u})^2, \quad u_i \in U \quad (2)$$

- 4: **While** the  $GVF$  is lower than maximum value **do**
- 5: Calculate the sum of squared deviation for each class,  $SD_j$ , using (3):
- 6: Increase the interval  $[INF_j, SUP_j]$  from classes with lowest  $SD_j$  by decreasing the interval from classes with largest  $SD_j$ .
- 7: Calculate the goodness of variance fit,  $GVF$ , using (4):

$$GVF = 1 - \frac{\sum_{j=1}^{NC} SD_j}{SD_U} \quad (4)$$

- 8: **End while**
- 9: Store the classes' boundaries of input dataset,  $U$ .

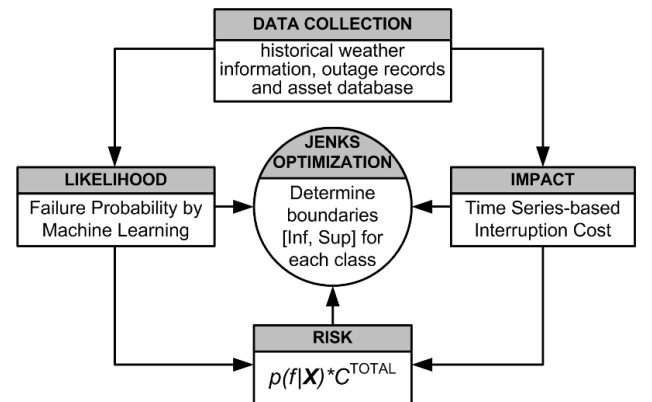


Fig. 1. Role of Jenks natural breaks optimization on risk assessment framework.

classification. The illustrated process to determine class boundaries can be a periodic procedure using, for instance, data collection from last year.

#### A. Failure Probability Metric by Machine Learning

The proposed risk assessment framework employs the failure probability metric to determine the likelihood of something is malfunctioning in a distribution network. The processing of large volume of data from diverse databases, i.e. outage management system (OMS), lighting detection network, GIS, weather stations, and asset management system (AMS) database, [15] and [16], contributes to threats characterization. Thus, the use of the big data analytics is required where the machine learning technique demonstrates great efficiency in the knowledge extraction. The Naive Bayes is the supervised learning technique used to establish an association of several features of interest into just one quantitative parameter [17]. In power distribution system, the failure probability metric uses a Naive Bayes model, as in (5) and (6), by taking into account the external dependences that are given by different types of threats as features of interest.

$$p(f | \mathbf{X}) = \frac{\tilde{p}(\mathbf{X} | f) \tilde{p}(f)}{\sum_f \tilde{p}(\mathbf{X} | f) \tilde{p}(f)} \quad (5)$$

$$\tilde{p}(\mathbf{X} | f) = \prod_{i=1}^D (\tilde{\theta}_{i,f}^{x_i} (1 - \tilde{\theta}_{i,f})^{1-x_i}) \quad (6)$$

where

- $p(f | \mathbf{X})$  Conditional probability of failure subject to  $\mathbf{X}$ ;
- $\tilde{p}(\mathbf{X} | f)$  Estimate of the likelihood of  $\mathbf{X}$  given  $f$ ;
- $\tilde{p}(f)$  Estimate of failure probability;
- $\tilde{\theta}_{i,f}^f$  Estimate of the probability of observing  $x_i$  conditioned to a failure event,  $f$ .

The two states of the failure feature,  $\text{dom}(f) = \{0,1\}$ , leads the definition of  $p(f = 1 | \mathbf{X})$  as the conditional probability of failure occurrence subject to observe external dependences,  $\mathbf{X} = \{x_i | \text{dom}(x_i) = \text{dom}(f) \wedge i = 1, \dots, D\}$ , that are enumerated in Table II. The probability of observing the vector  $\mathbf{X}$  can be compactly written as in (6) where  $p(x_i = 1 | f) \equiv \theta_{i,f}^f$  and  $p(x_i = 0 | f) \equiv 1 - \theta_{i,f}^f$  because of Naive Bayes conditional independence assumption. Another benefit of this assumption is the applying of maximum likelihood learning to the Naive Bayes model.

$$p(x_i = 1 | f) = \frac{\text{number times } x_i = 1 \text{ for } f}{\text{number of data points in } f} \quad (7)$$

TABLE II. OBSERVED EXTERNAL DEPENDENCES IN THE BAYES MODEL.

$x_i$	Feature of interest	$x_i$	Feature of interest
$x_1$	Wind speed is low	$x_6$	Weather is rainy
$x_2$	Wind speed is medium	$x_7$	Weather is thunderstorm
$x_3$	Wind speed is high	$x_8$	Incidence of lightning
$x_4$	Weather is good	$x_9$	Vegetation is over height
$x_5$	Weather is misty	$x_{10}$	Degradation by ageing

In the proposed model, the obtaining of the monthly likelihood to every features of interest, as in (7), determines the knowledge extraction from the available databases. Additionally, the prediction of the probability value in the current year of analysis is achieved using a regression model resulting of the ordinary least square (OLS) estimator, as given by (8) and (9).

$$\tilde{\theta}_{i,f}^f = \mathbf{B}_i^T \mathbf{T} \quad (8)$$

$$\mathbf{B}_i = \left( \sum_y \mathbf{T}_y (\mathbf{T}_y)^T \right)^{-1} \sum_y \mathbf{T}_y \theta_{i,y}^f \quad (9)$$

where

$\tilde{\theta}_{i,f}^f$  Estimate of the probability of observing  $x_i$  conditioned to a failure event,  $f$ , in the  $y^{\text{th}}$  year.

Since, the past years of observed probabilities are in the matrix  $\mathbf{T}_y = (1 \ t_y)^T$ , the prediction parameters comprise the matrix  $\mathbf{B}_i = (\beta_{0,i} \ \beta_{1,i})^T$ . The knowledge extraction from observed databases is achieved by calculating the prediction parameters in  $\mathbf{B}_i$ . The knowledge extraction is a function of data mining or knowledge discovery from data (KDD) that sequentially groups several functions for dealing with massive database difficulties, e.g. unnecessary information and inconsistent data [18]. In this way, the cleaning, integration and selection functions are performed before the knowledge extraction function that processes the useful information.

The processing of large volume of data also requires the integration of different sources of information. Fig. 2 shows the information flow for characterizing feeder sections of distribution networks in concordance with their failure probability. The distribution network operator workstation performs an important role by running the supervisory application, [19] and [20], with the addition of objects for regression and Bayes models.

Firstly, the estimative of probability values for each feature of interest uses the eq. (8). The elements of  $\mathbf{B}_i$  are obtained through the stored procedures in the historical database server. Secondly, the calculation of failure probability to every feeder sections is performed using eq. (5) and (6) where the vector of current external dependences, or observed statuses of features of interest, comes from external http servers for weather forecasting and lightning monitoring and from vulnerability models for vegetation and ageing.

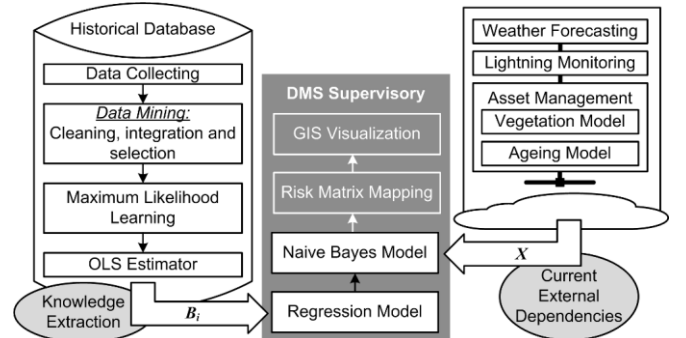


Fig. 2. Information flow involving the Naive Bayes probability estimation.

Many power flow interruptions that reduce the reliability indices are often caused when tree branches touch the distribution feeder conductors. The vegetation location detection is performed using remote sensing technology in association with GIS application that identifies the distribution feeder segments vulnerable to tree size. The prediction of tree heights uses a vegetation growth model as a function of time or age, [21] and [22], indicating whether computed tree height is over allowable height. Other vulnerability model takes into account electrical, mechanical and thermal stresses to determine the equipment degradation [23]. The ageing model makes use of the repair cycle for correlating equipment operating state and power supply interruption information. Thus, electrical equipment may have a high level of degradation whenever it reaches at least 63% of possibility to fail.

### B. Time Series-based Interruption Cost

In the proposed risk assessment approach the impact quantification is achieved by calculating the energy supply interruption cost [24]. The support of time varying energy consumption profiles is guaranteed by the time series-based interruption cost formulation as well as the identification of event locations involved in the outage management is supported by georeferenced network data. Considerable data on individual customers and power distribution system are required in the estimation of costs associated with the interruption. The utility company has costs that are related to income, electric energy sales, capital investments in their electrical devices and the operation and maintenance tasks. The regulatory authority maximizes the energy benefits to the society by balancing the energy consumption prices according to established rate-case rules. The energy purchase price and financial loss into power supply interruption also affects the customers' activities [25]. Hence, the sum of costs perceived by these various agents of the energy market yields the total cost of the power interruption,  $C^{\text{TOTAL}}$ , as in (10) - (14).

$$C^{\text{TOTAL}} = C^{\text{O\&M}}(\Delta d) + \sum_{K \in \Gamma} C^K(\Delta t) \quad (10)$$

$$C^K(\Delta t) = \sum_{i \in \Omega} \sum_{j \in \Phi} c_{i,j}^K z_{i,j} \quad (11)$$

$$c_{i,j}^{\text{ens}} = c_j^e L_j \sum_{m \in \Theta} \sum_{n \in T} F_{i,m,n}^{\text{dem}} w_{j,m,n} \quad (12)$$

$$c_{i,j}^{\text{pen}} = H(i\partial t, \Delta t^{\text{max}}) c_{i,j}^{\text{ens}} \quad (13)$$

$$c_{i,j}^{\text{ic}} = L_j \sum_{m \in \Theta} \sum_{n \in T} (c_{i,m}^{\text{CDF}} - c_{i-1,m}^{\text{CDF}}) F_{i,m,n}^{\text{dem}} w_{j,m,n} \quad (14)$$

A feeder section comprises of several electrical components permitting the network reconfiguration by opening and closing emergency connectors, or sectionalizing switches, in the terminations connected with other feeder sections. In (10), the total cost caused by the interruption of one feeder section is given in two parts. The first part is the operation and maintenance cost,  $C^{\text{O\&M}}$ , that depends on the route traveled by

the field crew,  $\Delta d$ , where the distribution network topology, georeferenced position of sectionalizing switches, initial position of field crews and GIS routing application are employed as input information for solving the crew dispatch problem [26]. The second part is the sum of cost related to different market agents that are grouped in a set  $\Gamma = \{\text{ens}, \text{pen}, \text{ic}\}$  comprising, respectively, the billing loss of utility company, the penalty cost from regulatory authority rules, and the economic losses of different types of customers.

These different costs,  $C^K$ , depend on the interruption time,  $\Delta t$ , i.e. the time span including outage report time (wait time from the fault occurrence until the dispatch of field crews), maneuver time (interval involving the field crew travel, feeder inspection and manual switching to isolate the faulted feeder section and to restore the adjacent feeder sections) and repair time (required time to repair the damage equipment and to restore the energy supply service). Since fault management procedures change the state of energy customers, the interruption time is discretized by a pre-defined time step,  $\partial t$ , yielding the set of time series,  $\Omega$ . In (11),  $z_{i,j}$  is a binary variable that reproduces state changes of the  $j^{\text{th}}$  customer during the interruption time where the logic value 1 indicates the energy supply interruption. The  $\Phi$  set contains all customers on the feeder and the effect of different market agents over the individual customer cost,  $c_{i,j}^K$ , follows the formulation as given in (12) - (14).

Additionally to operation and maintenance cost, the utility company also perceives the billing loss, i.e. the cost of energy that could be sold to customers during the interruption, given by the cost of energy not supplied,  $c_{i,j}^{\text{ens}}$ , as (12) where  $c_j^e$  is electricity rate and  $L_j$  is the installed power of the  $j^{\text{th}}$  customer. The most typical customer types are grouped in  $\Theta = \{\text{residential}, \text{commercial}, \text{industrial}\}$  while their consumption profiles are in  $T = \{\text{low}, \text{medium}, \text{high}\}$ . In this way,  $F_{i,m,n}^{\text{dem}}$  is a tridimensional data array with load percentage demand hour-by-hour [24] and, consequently,  $w_{j,m,n}$  is a two-dimensional binary array for indicating the type and consumption profile of the  $j^{\text{th}}$  customer.

According to the rules established by regulatory authorities for compensating customers over long outages [27], utility companies could be penalized and customer compensated whenever the outage interval exceeds the established limit. In (13), the penalty cost,  $c_{i,j}^{\text{pen}}$ , is determined using the  $H$  function that has zero value while the product of  $i\partial t$  is less than the maximum outage duration,  $\Delta t^{\text{max}}$ . Otherwise, the billing loss of  $j^{\text{th}}$  customer is multiplied by a factor of penalty.

The most significant part of the total cost is the customer interruption cost that associates the economic losses of different customers during the power supply failures [28]. Wages paid to idle workers, loss of sales, overtime costs, damage to equipment, spoilage of perishables, cost of running back-up generators and cost of any special business procedures contribute to the determination of the customer interruption cost [29]. In particular, the endangered well-being, spoiled food and damaged appliances may affect

residential customers. The impact of power interruption is popular and directly formulated using the customer damage function by expressing the customer interruption cost as a function of outage duration [30]. Equation (14) determines the customer interruption cost,  $c_{i,j}^{ic}$ , for  $j^{th}$  customer in the  $i^{th}$  time step. The values of  $c_{i,m}^{CDF}$  time series are interpolations from the table containing values of customer damage functions that are typically defined for each economic activity or customer type.

### C. Method for Defining the Risk Matrix

The calculation of failure probability and interruption cost quantifies the likelihood and impact, respectively, and should be performed hour-by-hour for timely risk assessment using the proposed risk matrix. Hourly values of likelihood and impact are classified according to the categories determined by Jenks optimizer and mapped to rows and columns of the risk matrix whose elements determine risk levels. Hence, assigning the risk level for each element in the risk matrix is fundamental to the risk assessment effectiveness. This process can be updated annually using information from last year to update the risk matrix for current year.

As demonstrated in Fig. 1, the risk is also quantified by multiplying  $p(f|X)$  times  $C^{TOTAL}$  and classified in risk levels using the Jenks natural breaks algorithm. If quantified values of likelihood and impact from previous year are disposed into a dispersion chart, the result can be presented in form of the graphic where each data point is classified according to risk levels. Since likelihood and impact categories are limited and cover axes of dispersion chart, there is a limited number of discrete regions as, for example, the region that is limited by categories II and B. This region determines the value of an element into risk matrix because its rows and columns are mapped by likelihood and impact categories; however, some regions in the dispersion chart have data points with different classification, for example, data points in a particular region can be classified as both medium and high risk level.

$$r_{m,n} = i, \text{ if } \rho_i = \max_{j=\{L,M,H\}} \left\{ \rho_j = \frac{\sum_{k \in \Pi(m,n)} dp_k |_{\text{level}=j}}{\sum_{k \in \Pi(m,n)} dp_k} \right\} \quad (15)$$

The element,  $r_{m,n}$ , of the risk matrix is then determined using the density method as is formulated in (15) where the value of  $i$  is equal to the risk level ( $L$ ,  $M$  or  $H$ ) with the maximum calculated density,  $\rho_i$ , in the  $\Pi(m,n)$  region that is limited by  $m^{th}$  likelihood category and  $n^{th}$  impact category. In other words,  $r_{m,n} = L$  if the number of data points classified as low risk level,  $dp_k |_{\text{level}=L}$ , is preponderant in the region  $\Pi(m,n)$ . In the region with identical values of calculated densities, the value of the element representing this region into risk matrix is equal to the highest risk level because higher risk levels are less frequent than lower risk levels.

The determination of risk matrix elements completes the inference mechanism of the proposed online risk assessment for each feeder section of power distribution network. Although formulated models are very important in the

quantification of likelihood and impact, the central issue in this work relates to process of how to classify these quantities, how to build the risk matrix and how to develop a DMS tool able to efficiently display the risk levels using a GIS application. Therefore, the following section comprises both the construction of risk matrix by determining classes' boundaries and the verification of the developed GIS tool for risk assessment.

## IV. GIS VISUALIZATION BY RISK MATRIX

The proposed methodology is evaluated under real world distribution feeder with data available in [31]. Ten sectionalizing switches limits nine feeder sections in the evaluated feeder. These feeder sections have multiple laterals and electrical loads and are also limited by sectionalizing switches that must operate during the reconfiguration procedure. In the calculation of failure probability, the learning information comes from external sources: two weather stations and one lightning detection network, where the historical databases comprise seven years, from 2009 to 2015. Parameters of the vegetation growth model are adjusted by considering the tree pruning schedule equals to one year whereas the equipment degradation vulnerability model of different devices may have their parameters obtained using the method discussed in [23]. In terms of interruption cost, the input dataset can be found in [24]. Both calculations obtain quantified values of likelihood and impact for each feeder section. A general purpose programming language (C++) is used in the implementation of the proposed models that are integrated with a distribution network simulation platform for supporting the use georeferenced data [20].

### A. Building the Risk Matrix

Fig. 1 illustrates a logic diagram with processes for building the risk matrix using last year's collected data. *The first process* comprises the determination of quantified likelihood ranges by defining inferior and superior boundaries through Jenks optimization. The goodness of variance fit (GVF) is a quality index used by the Jenks algorithm as stopping criteria. The perfect fit, or "optimum data classification", is achieved when  $GVF = 1$ . In the classification process, the histogram was built using around five thousand values of failure probability.

Fig. 3 shows the histogram of the distribution of failure probabilities where the frequency axis is rated using logarithmic scale of base ten. A histogram in linear scale is shown at the far-right corner, which helps to deduce the absence of a probability density function able to characterize the likelihood. There are failure probability values with zero frequency because the set of external dependences,  $X$ , has a finite number of features of interest and the occurrence probability for each feature of interest is calculated monthly. Despite this characteristic, the Jenks optimizer found the six likelihood categories and their range limits by a GVF index being equal to 0.98704. For example, the likelihood category *III* comprises failure probability values between 0.31 and 0.54.

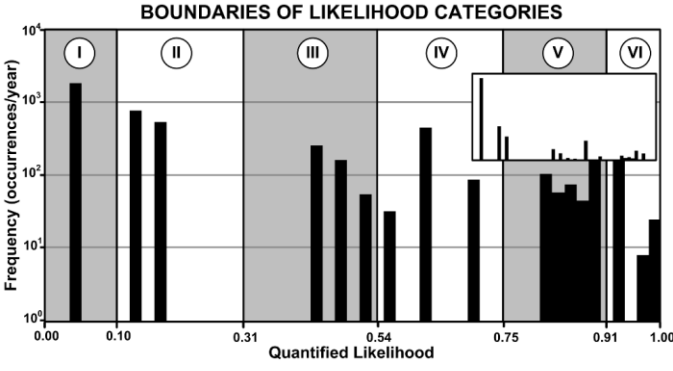


Fig. 3. Graphical representation of the distribution of failure probability values with likelihood category ranges.

The second process involves the determination of quantified impact categories by determining their boundaries. Fig. 4 displays the histogram of the distribution of interruption cost values where frequency was obtained by taking into account a series of intervals each equal to \$500. The distribution characteristic is shown by the histogram in linear scale helping to deduce that interruption cost values can be featured by a Weibull probability distribution. Although the economic activity and consumption profile are important factors in the cost calculation, the interruption duration, which also figures the Weibull probability density function, is the factor with the greatest influence over the interruption cost.

Six impact categories were achieved by Jenks algorithm with GVF index equals to 0.96149. The first three impact categories have shorter ranges due to large frequencies in this region. Consequently, the impact category B has the shortest range, equals to \$3500, whereas the category F comprise the longest range, from \$20,450 to \$31,670.

Once the risk is quantified by multiplying failure probabilities times interruption costs, the third process deals with the determination of risk levels by defining their boundaries. Fig. 5 demonstrates the histogram of the distribution of quantified risk values using a series of intervals equal to \$500. The linear scale-based histogram at far-end right corner reveals that risk distribution has the behavior of an exponential probability distribution, so the most adequate classification methodology should be performed by head/tail breaks classifier [32]. In this case, the Jenks optimizer can be used again because the quantified risk is grouped in few numbers of classes, i.e. in three risk levels, and the density

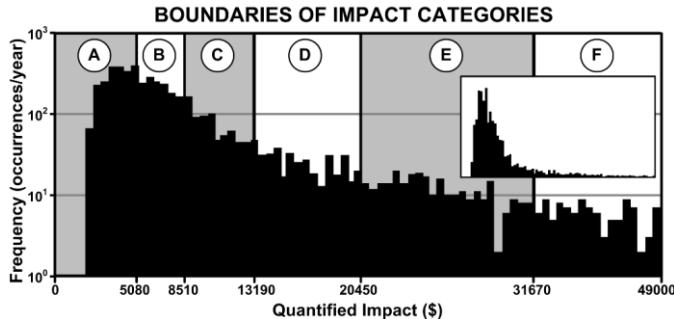


Fig. 4. Graphical representation of the distribution of interruption cost values with impact category ranges.

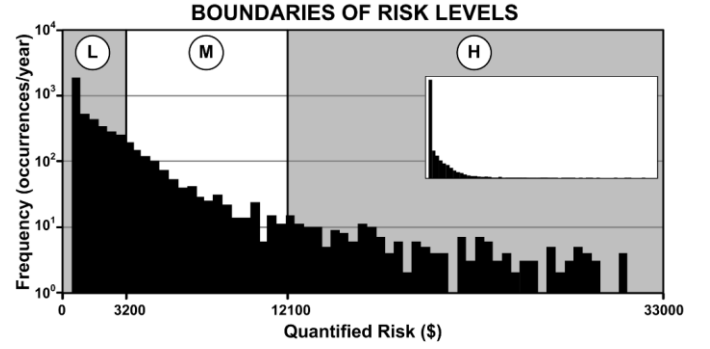


Fig. 5. Graphical representation of the distribution of quantified risk values with risk level ranges.

method should still determine the preponderant characteristic for each region at dispersion chart what admits data points that are classified with less degree of accuracy.

Three risk levels were achieved using the Jenks algorithm with GVF index equals to 0.83967. Although the quality index had been worse than GVF indices in quantified likelihood and impact classification, the achieved risk level ranges fit with heavy-tailed distribution. For instance, the head risk level, L, has range equals to \$3,200 in contrast to the tail risk level, H, with range of \$21,900.

After the determination of class boundaries, the next process consists of the construction of risk matrix using the density method. Table III presents elements of the risk matrix where rows are likelihood categories and columns are impact categories. Now, the hourly risk assessment can be executed using previously determined categories and risk matrix.

#### B. Study Case under Real Distribution Network

Fig. 2 presents displays that the distribution operator will see at DMS supervisory running the GIS web application with risk matrix mapping. The developed GIS application performs likelihood and impact quantification and classification for each feeder section of power distribution network. Two achieved categories define one row and one column in the risk matrix whose intersection determines the risk level of the feeder section.

According to Table III, each risk level is identified by a color, thus, the GIS application assigns for the graphical representation of the feeder section the color corresponding to the risk level. Furthermore, the addition of daily hours to set of spatial coordinates includes one more dimension into feeder section representation in GIS application. This extra dimension has the risk level information represented hour-by-hour, which is well suited to perform online risk mitigation.

TABLE III. DETERMINED ELEMENTS OF THE RISK MATRIX.

		IMPACT					
		A	B	C	D	E	F
LIKELIHOOD	I	L	L	L	L	L	L
	II	L	L	L	L	L	M
	III	L	L	M	M	M	H
	IV	L	M	M	M	H	H
	V	M	M	M	H	H	H
	VI	M	M	M	H	H	H



Fig. 6. Partial screen of the developed GIS application with tridimensional representation of risk levels hour-by-hour.

Fig. 6 shows a screen shot with the tridimensional graphical representation of the tested distribution network where different colors are hourly risk levels. The base of the graphic corresponds to daily early hours, from 00:00 to 06:00 of January, 14th of 2016, with low risk level in all feeder sections. After that, both weather condition and energy consumption profile are modified causing changes in risk level of feeder sections. For example, feeder section #1 presents very low risk level but, along the day, its risk level was classified as medium because of weather changes. At 20:00, the observed weather pattern was thunderstorm with medium wind speed given by  $X = \{0100001x_8x_9x_{10}\}$  causing the feeder section #2 to change its risk level from medium to high risk. Although weather changes influence the risk level in feeder section #3, the main color is intense red representing the high risk level that is a consequence of economic activities from customers with large installed power.

The other way of taking advantage of the developed GIS tools is by assigning the value attribution to risk levels, for instance, low level is equals to 0, medium is equals to 1 and high is 2. Thus, the different grades of the accumulated risk along the distribution network are visualized using color

temperature scale in overlapped layers with different accumulated risk values. Fig 7 shows the accumulated risk values during January where the smaller accumulated risk values are the first layers in cold color while the larger values are the last layers in hot color. The feeder section #1 has one lower layer in cold color indicating the accumulated risk is small. On the other hand, the feeder section #2 had upper layers with hot color tones indicating its large accumulated risk, which is the consequence of customers' types connected in this section.

The high risk level does not just depend on the failure probability but also on the impact intensity, as is established in Table III. When the failure probability quantization has large value and it is classified as Likely (VI), the risk level must be either medium (M) or high (H). In the comparison process, the existence of low (L) risk level at a failure event indicates hence a mismatching of the proposed methodology. Fig. 8 shows that the proposed methodology presents a mismatching ratio around 20% whenever the cause of failure event is adverse weather, component failure or lightning. When the cause is vegetation contact, the ratio improves to 10%. Subsequently, the hours after one mismatching the correct risk

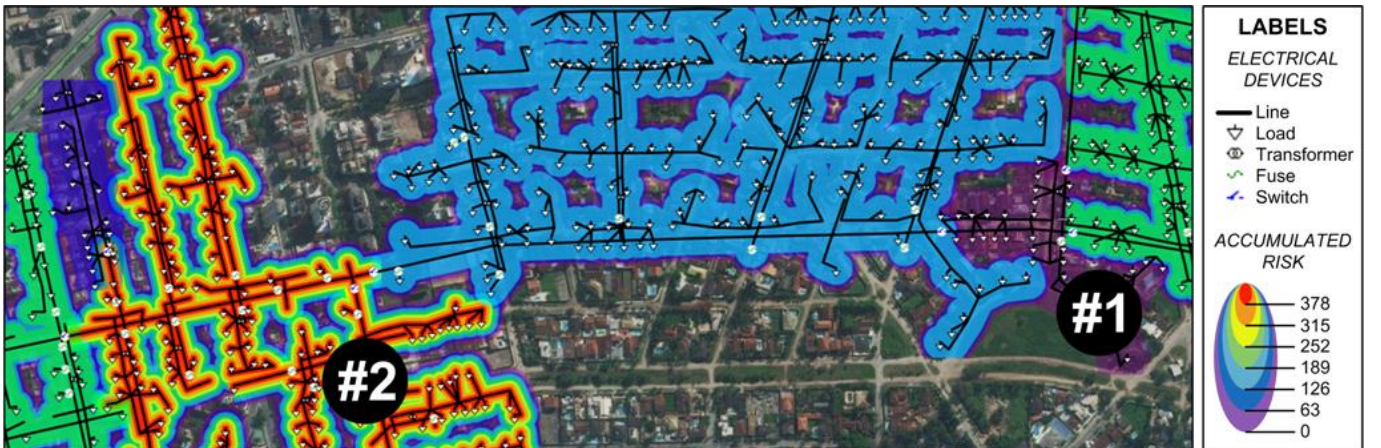


Fig. 7. Partial screen of the developed GIS application with tridimensional representation of accumulate risk levels during a month.

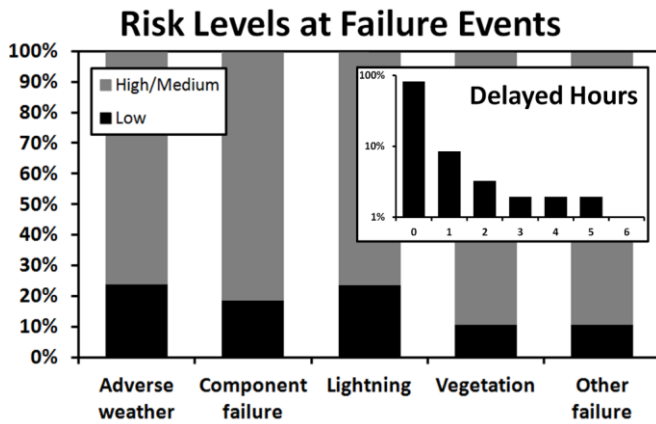


Fig. 8. Graphical representation of the comparison process with the percentages of risk levels at failure events.

level are calculated and indicated by the proposed methodology. The bar chart of delayed hours demonstrates that the delay time does not overcome five hours and in the most part of mismatching occurrences the correct risk level is indicated with one hour of delay. These results reveal the effectiveness of the proposed methodology for evaluating the operating condition of power distribution networks.

## V. CONCLUSION

We have shown that the weather-based risk assessment can provide risk quantification through the correlation involving available weather data and historical management data of the power distribution system.

Once the realization of this risk assessment step is implemented, one can then integrate it with the advanced distribution management system to offer risk mitigation. This tool facilitates the operators' decisions since it employs spatiotemporal GIS based visualization of the resiliency improvement actions.

## REFERENCES

- [1] D. Vose, *Risk Analysis: a Quantitative Guide*, 3rd ed., vol. 1, West Sussex, England: John Wiley & Sons, 2008, 729 pp.
- [2] T. Dokic, P. Dehghanian, P. Chen, M. Kezunovic, Z. Medina-Centina, J. Stojanovic and Z. Obradovic, "Risk Assessment of a Transmission Line Insulation Breakdown due to Lightning and Severe Weather," in *Proc. HICSS*, Kauai, HI, 2016, pp. 2488-2497.
- [3] S. Nattian and M. Kezunovic, "A Risk-Based Approach for Maintenance Scheduling Strategies for Transmission System Equipment," in *Proc. PMAPS*, Rincón, Puerto Rico, 2008, pp. 1-6.
- [4] R. Billinton and Z. Pan, "Historic Performance-based Distribution System Risk Assessment," *IEEE Trans. Power Del.*, vol. 19, no. 4, pp. 1759-1765, Oct. 2004.
- [5] J. L. Pinto and M. Louro, "On Human Life Risk-Assessment and Sensitive Ground Fault Protection in MV Distribution Networks," *IEEE Trans. Power Del.*, vol. 25, no. 4, pp. 2319-2327, Oct. 2010.
- [6] T. V. Garcez and A. T. de Almeida, "Multidimensional Risk Assessment of Manhole Events as a Decision Tool for Ranking the Vaults of an Underground Electricity Distribution System," *IEEE Trans. Power Del.*, vol. 29, no. 2, pp. 624-632, Apr. 2014.
- [7] Y. Ding, M. Xie, Q. Wu and J. Ostergaard, "Development of Energy and Reserve Pre-Dispatch and Re-Dispatch Models for Real-Time Price Risk and Reliability Assessment," *IET Gener. Transm. Distrib.*, vol. 8, no. 7, pp. 1338-1345, Jul. 2014.
- [8] R. Shah, R. Yan and T. K. Saha, "Chronological Risk Assessment Approach of Distribution System with Concentrated Solar Power Plant," *IET Gener. Transm. Distrib.*, vol. 9, no. 6, pp. 629-637, Aug. 2015.
- [9] W. Deng, H. Ding, B. Zhang, X. N. Lin, P. Bie and J. Wu, "Multi-Period Probabilistic-Scenario Risk Assessment of Power System in Wind

- Power Uncertain Environment," *IET Gener. Transm. Distrib.*, vol. 10, no. 2, pp. 359-365, Feb. 2016.
- [10] M. Nijhuis, M. Gibescu and S. Cobben, "Risk-based Framework for the Planning of Low-Voltage Networks Incorporating Severe Uncertain," *IET Gener. Transm. Distrib.*, vol. 11, no. 2, pp. 419-426, Jan. 2017.
- [11] B. M. Ayyub, "Risk Analysis Methods," in *Risk Analysis in Engineering and Economics*, 1st ed., vol. 1, Boca Raton, FL: Chapman & Hall/CRC, 2003, pp. 33-118.
- [12] G. F. Jenks, "The Data Model Concept in Statistical Mapping," in *International Yearbook of Cartography*, vol. 7, Liverpool, England: George Philip & Son Ltd., 1967, pp. 186-190.
- [13] T. Slocum, R. McMaster, F. Kessler and H. Howard, *Thematic Cartography and Geovisualization*, 3rd ed., vol. 1, Upper Saddle River, NJ: Pearson Prentice Hall, 2009, p. 561.
- [14] J. B. MacQueen, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, Berkeley, CA: Univ. of Calif. Press, 1967, pp. 281-297.
- [15] P. Chen, T. Dokic and M. Kezunovic, "The Use of Big Data for Outage Management in Distribution System," in *Proc. CIGRE*, Rome, Italy, 2014, pp. 1-5.
- [16] M. Kezunovic, L. Xie and S. Grijalva, "The Role of Big Data in Improving Power System Operation and Protection," in *Proc. IREP*, Rethymnon, Greece, 2013, pp. 1-9.
- [17] D. Lowd and P. Domingos, "Naive Bayes Models for Probability Estimation," in *Proc. ICML*, Bonn, Germany, 2005, pp. 1-8.
- [18] J. Han, M. Kamber and J. Pei, *Data Mining: Concepts and Techniques*, 3rd ed., vol. 1, Waltham, MA: Elsevier Inc., 2012, p. 740.
- [19] J. B. Leite and J. R. S. Mantovani, "Development of a Smart Grid Simulation Environment, Part I: Project of the Electrical Devices Simulator," *J. Control Autom. Electr. Syst.*, vol. 26, no. 1, pp. 80-95, Feb. 2015.
- [20] J. B. Leite and J. R. S. Mantovani, "Development of a Smart Grid Simulation Environment, Part II: Implementation of the Advanced Distribution Management System," *J. Control Autom. Electr. Syst.*, vol. 26, no. 1, pp. 96-104, Feb. 2015.
- [21] D. T. Radmer, P. A. Kuntz, R. D. Christie, S. S. Venkata and R. H. Fletcher, "Predicting Vegetation-Related Failure Rates for Overhead Distribution Feeders," *IEEE Trans. Power Deliv.*, vol. 17, no. 4, pp. 1170-1175, Oct. 2002.
- [22] F. B. Martins, C. P. B. Soares and G. F. da Silva, "Individual Tree Growth Models for Eucalyptus in Northern Brazil," *Sci. Agric.*, vol. 71, no. 3, pp. 212-225, May 2014.
- [23] X. Zhang and E. Gockenbach, "Component Reliability Modeling of Distribution Systems based on Evaluation of Failure Statistics," *IEEE Trans. Dielec. And Elect. Insul.*, vol. 14, no. 5, pp. 1183-1191, Oct. 2007.
- [24] J. B. Leite, J. R. S. Mantovani, T. Dokic, Q. Yan, P. -C. Chen and M. Kezunovic, "The Impact of Time Series-Based Interruption Cost on Online Risk Assessment in Distribution Networks," in *Proc. T&D LA*, Morelia, Mexico, 2016, pp. 1-6.
- [25] M. J. Sullivan, B. N. Suddeth, T. Vardell and A. Vojdani, "Interruption Costs, Customer Satisfaction and Expectations for Service Reliability," *IEEE Trans. Power Syst.*, vol. 11, no. 2, May 1996.
- [26] P. M. S. Carvalho, F. J. D. Carvalho and L. A. F. M. Ferreira, "Dynamic Restoration of large-Scale Distribution Network Contingencies: Crew Dispatch Assessment," in *Proc. PowerTech*, Lausanne, Switzerland, 2007, pp. 1453-1457.
- [27] C. J. Wallnerstrom and P. Hilber, "Vulnerability Analysis of Power Distribution Systems for Cost-Effective Resource Allocation," *IEEE Trans. Power Syst.*, vol. 27, no. 1, pp. 224-232, Feb. 2012.
- [28] Q. Yan, T. Dokic and M. Kezunovic, "Predicting Impact of Weather Caused Blackouts on Electricity Customers Based on Risk Assessment," in *Proc. PESGM*, Boston, MA, 2016, pp. 1-6.
- [29] N. Kaur, G. Singh, M. S. Bedi and T. Bhatti, "Evaluation of Customer Interruption Cost for Reliability Planning of Power System in Developing Economies," in *Proc. PMAPS*, Ames, IA, 2004, pp. 752-755.
- [30] R. Billinton and W. Wangdee, "Approximate Methods for Event-Based Customer Interruption Cost Evaluation," *IEEE Trans. Power Syst.*, vol. 20, no. 2, pp. 1103-1110, May 2005.
- [31] S. Commission. (2016, Mar.). Distribution testing system of 1807 lines. UNESP, Ilha Solteira, Brazil. [Online]. Available: [http://www.feis.unesp.br/Home/departamentos/engenhariaeletrica/lapsee/807/home/distribution\\_network\\_1806\\_lines.rar](http://www.feis.unesp.br/Home/departamentos/engenhariaeletrica/lapsee/807/home/distribution_network_1806_lines.rar).
- [32] B. Jiang, "Head/Tail Breaks: A New Classification Scheme for Data with a Heavy-Tailed Distribution," *J. Prof. Geogr.*, vol. 65, no. 3, pp. 482-494, Jul. 2013.