





Deciphering the drivers of direct and indirect damages to companies from an unprecedented flood event: A data-driven, multivariate probabilistic approach

Ravi Kumar Guntu^{1*}, Guilherme Samprogna Mohor², Annegret H. Thieken², Meike Müller³, and Heidi Kreibich¹

¹GFZ Helmholtz Centre for Geosciences, Section Hydrology, 14473 Potsdam, Germany

²Institute of Environmental Science and Geography, University of Potsdam, 14476 Potsdam, Germany

³Deutsche Rückversicherung AG, 40549 Düsseldorf, Germany

*Correspondence to: Ravi Kumar Guntu (ravikumar.guntu@gfz.de)

1.1 Precaution ratio

Particular	Survey Question	If	Conversion	Else
v1	Has the heating system been adjusted?	implemented before the event in July 2021	1	0
v2	Is there an emergency plan as of July 2021?	Yes	1	0
v3	What was the frequency of drills conducted before July 2021?	less than once a year	1	
		once a year	1	
		twice a year	1	
		more than twice a year	1	
v4	Are tanks etc. anchored securely?	implemented before the event in July 2021	1	0
v5	Stationary or mobile water barriers?	implemented before the event in July 2021	1	0
v6	Sensitive equipment on higher floors?	implemented before the event in July 2021	1	0
v7	Water-hazardous substances on higher floors?	implemented before the event in July 2021	1	0
v8	Area used adjusted to risk?	implemented before the event in July 2021	1	0
v9	Is the air conditioning equipped for flood safety?	implemented before the event in July 2021	1	0
v10	Has flood safety been improved?	implemented before the event in July 2021	1	0

Weighted Sum Calculation

The weighted sum (p) is calculated as:

$$p = v1 + v2 + v3 + (5 \times v4) + (5 \times v5) + (5 \times v6) + (5 \times v7) + (10 \times v8) + (10 \times v9) + (10 \times v10)$$

Precaution Indicator Based on Weighted Sum

Precaution Indicator	Weighted Sum (p)
0 (No)	0
1 (Medium)	Between 1 and 5
2 (Good)	6 or greater

Damage type	bdam(188)	0.04	0.02	0.01	0.01	0.1	0.05	0.01	0	0.01	0.01	0	0.01	0	0.02	0.03	0	0	0.01	0.02
	edam(272)	0.02	0.04	0.03	0.01	0.1	0.04	0	0	0	0.04	0	0.03	0	0.04	0.05	0	0	0	0
	gsdam(217)	0.02	0.06	0.02	0.01	0.09	0.04	0	0	0	0.04	0	0.03	0	0.02	0.03	0	0	0	0
	bid(332)	0.03	0.07	0.04	0.03	0.12	0.04	0	0	0.01	0.04	0	0.03	0	0.06	0.07	0	0	0.02	0.02
	brd(257)	0.04	0.07	0.04	0.03	0.1	0.04	0	0	0.01	0.04	0	0.04	0	0.07	0.07	0	0	0.01	0.01
		wd	d	v	con	wt	ws	ew	me	ep	kh	ms	fe	pr	in	sp	sec	ss	own	emp
		Variable																		

Figure S1: Percentage of missing values per factor (x-axis) for each damage type (y-axis). The value in parentheses for each damage type indicates the number of responses available out of 431. For warning time (wt), cases where no warning was received are treated as zero.

1.2 Optimal hyperparameters selection

We conducted hyperparameter tuning, which randomly samples from a predefined range of values for each hyperparameter. The search was performed within a nested cross-validation framework. The framework consists of two levels of cross-validation.

Outer cross-validation: The outer loop is responsible for evaluating the overall performance of the model. It splits the dataset into training and testing subsets, ensuring that performance estimation is conducted on unseen data. We employed *RepeatedKfold* cross-validation with 10 splits and 10 repeats, resulting in a total of 100 evaluations.

Inner cross-validation: Within each outer loop iteration, the inner cross-validation loop is used to tune hyperparameters. We employed *RandomizedSearchCV*, which samples different hyperparameter combinations and evaluates them using 10-fold cross-validation within the training set.

This approach ensures that model selection (inner loop) and performance evaluation (outer loop) are conducted independently, leading to a more reliable assessment of the model's generalization ability (Cawley and Talbot, 2010). The optimal combination of hyperparameters was selected based on the lowest Mean Absolute Error (MAE) during outer cross-validation.

Example breakdown for one outer fold – bdam (total 188 samples)

Outer loop (10-fold cross-validation, repeated 10 times)

- Test set: 19 samples (1-fold, unseen during training and tuning).
- Training set: 169 samples (remaining 9 folds, used for inner-loop tuning and training).

Inner loop (10-fold cross-validation within the training set of 169 samples)

- Training fold (inner loop): 152 samples (9 folds, used for model training).
- Validation fold (inner loop): 17 samples (1-fold, used to evaluate hyperparameters tuning).

The following hyperparameters and their respective search ranges were examined. All other hyperparameters were set to their default values as provided in the relevant scikit-learn packages (Pedregosa et al., 2011). The random state/seed was kept constant across all three models.

1.2.1 Elastic Net (EN) Hyperparameters

Elastic Net is a linear model that combines Lasso and Ridge regularization. The following hyperparameters control its regularization strength:

- α : Controls the strength of the regularization and balances between Lasso ($\alpha = 1$) and Ridge ($\alpha = 0$) regression. α : [0.01, 0.1, 0.25, 0.5, 0.75, 0.9, 1]
- λ : Determines the balance between the Lasso and Ridge terms. λ : [0.01, 0.1, 0.3, 0.5, 0.7, 0.9, 1.0]

1.2.2 Random Forest (RF) Hyperparameters

Random Forest is an ensemble method that constructs multiple decision trees. Key hyperparameters that influence the model's performance are:

- `n_estimators`: The number of trees in the forest. `n_estimators` = [50, 100, 200, 300]
- `max_depth`: The maximum depth of the individual trees. `max_depth` = [3, 5, 7, 10, 15, 20]
- `min_samples_split`: Minimum number of samples required to split an internal node. `min_samples_split` = [2, 5, 10]
- `min_samples_leaf`: Minimum number of samples required to be at a leaf node. `min_samples_leaf` = [1, 2, 4]
- `max_features`: Number of features to consider when splitting a node. `max_features` = ["sqrt", "log2", 0.5, 0.75, 1.0]

1.2.3 XGBoost (XGB) Hyperparameters

XGBoost is a gradient boosting method that optimizes speed and performance. The following hyperparameters are crucial for controlling the boosting process:

- `n_estimators`: The number of boosting rounds or trees. `n_estimators` = [50, 100, 200, 300, 500]
- `learning_rate`: Step size for each iteration, which controls the contribution of each tree. `learning_rate` = [0.005, 0.01, 0.05, 0.1, 0.2]
- `max_depth`: The maximum depth of each tree. `max_depth` = [3, 5, 7, 10, 15, 20]
- `subsample`: Fraction of samples used to train each tree. `subsample` = [0.1, 0.3, 0.5, 0.7, 1.0]
- `colsample_bytree`: Fraction of features used to train each tree. `colsample_bytree` = [0.1, 0.3, 0.5, 0.7, 1.0]
- `gamma`: Regularization parameter that controls overfitting. `gamma` = [0, 0.1, 0.5, 1, 5]
- `lambda` (L2 regularization) and `alpha` (L1 regularization): These terms control overfitting by penalizing the tree's complexity. λ = [0, 0.1, 1, 10], α = [0, 0.1, 1, 10].

1.3 Bayesian Network learning using the BDe score

The top four variables identified from the combined weighted importance scores obtained using three machine learning models were selected to predict flood damage through Bayesian Networks (BNs). Given that the dataset includes both discrete and continuous variables, continuous variables were discretized using an equal-frequency binning approach, ensuring all variables became categorical after discretization.

1.3.1 Bayesian Network Structure Learning

The Tabu search algorithm was applied to determine the Bayesian Network structure by maximizing the Bayesian Dirichlet equivalent (BDe) score. The learning process began with a simple Naive Bayes structure, in which all predictor variables (parent nodes) are assumed to be conditionally independent of each other given the target variable (child node). This initial structure reflects the assumption that no dependencies exist among predictors when conditioned on the target.

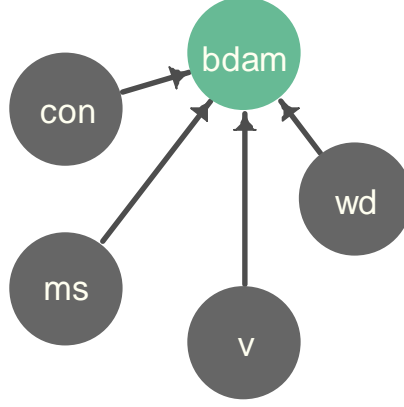


Figure S2: Illustration of the initial Naive Bayes structure for modeling company building damage using predictor variables as parents and the damage variable as the child node.

1.3.2 BDe Score Definition

The BDe score, a log-likelihood-based metric, evaluates how well a network structure (DAG) fits the data. It is computed as (Heckerman et al., 1995):

$$BDe\left(\frac{G}{D}\right) = \prod_{i=1}^n \left(\prod_{j=1}^{q_i} \left(\frac{\Gamma(\alpha_{ij})}{\Gamma(\alpha_{ij} + N_{ij})} \prod_{k=1}^{r_i} \frac{\Gamma(\alpha_{ijk} + N_{ijk})}{\Gamma(\alpha_{ijk})} \right) \right)$$

Where:

- n = number of variables
- q_i = number of parent configurations for variable X_i
- r_i = number of bins for X_i
- N_{ijk} = observed count of instances where $X_i = x_k$ given parent configuration $P a_j$
- N_{ij} = total count for parent configuration j .
- α_{ijk} = Dirichlet prior parameter.
- α_{ij} = Sum of Dirichlet priors for a given parent configuration.
- $\Gamma(\cdot)$ = Gamma function.

A higher BDe score indicates a better-fitting network.

Example calculation of the BDe score

Consider a simplified BN with two variables: *wd* (predictor: water depth) and *bdam* (target: relative damage to building). Both are discretized into bins as shown below.

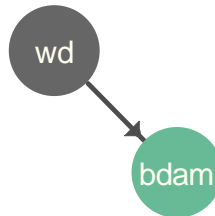


Figure S3: A simple Bayesian network representing the relationship between water depth (*wd*) and building damage (*bdam*).

1.3.3 Frequency Counts and Total Counts

Observed frequencies (N_{ijk}) from the sample data are presented in Table S1.

Table S1: Observed frequency of *bdam* conditional on *wd*.

<i>wd</i>	<i>bdam</i>	Count (N_{ijk})
$\leq 1\ m$	[0.00025,0.120)	30
$\leq 1\ m$	[0.12000,0.300)	19
$\leq 1\ m$	[0.30000,0.528)	7
$\leq 1\ m$	[0.52761,0.960)	4
$\leq 1\ m$	[0.96000,1.000]	5
$> 1\ m \ \& \ \leq 2\ m$	[0.00025,0.120)	3
$> 1\ m \ \& \ \leq 2\ m$	[0.12000,0.300)	18
$> 1\ m \ \& \ \leq 2\ m$	[0.30000,0.528)	16
$> 1\ m \ \& \ \leq 2\ m$	[0.52761,0.960)	20
$> 1\ m \ \& \ \leq 2\ m$	[0.96000,1.000]	16
$> 2\ m$	[0.00025,0.120)	5
$> 2\ m$	[0.12000,0.300)	4
$> 2\ m$	[0.30000,0.528)	18
$> 2\ m$	[0.52761,0.960)	7
$> 2\ m$	[0.96000,1.000]	16

Table S2: Total counts (N_{ij}) for each parent configuration (*wd* bins).

<i>wd</i>	Total N_{ij}
$\leq 1\ m$	65
$> 1\ m \ \& \ \leq 2\ m$	73
$> 2\ m$	50

1.3.4 Dirichlet Prior

Assuming a uniform Dirichlet prior with $\alpha = 15$:

$$\alpha_{ijk} = \frac{\alpha}{r \times q} = \frac{15}{5 \times 3} = 1$$

$$\alpha_{ij} = \sum_{k=1}^r \alpha_{ijk} = 5$$

1.3.5 BDe Score Computation

the BDe score is calculated for each parent configuration.

$$\text{For } wd \ (j = 1): \frac{\Gamma(5)}{\Gamma(5+65)} \times \frac{\Gamma(1+30)}{\Gamma(1)} \times \frac{\Gamma(1+19)}{\Gamma(1)} \times \frac{\Gamma(1+7)}{\Gamma(1)} \times \frac{\Gamma(1+4)}{\Gamma(1)} \times \frac{\Gamma(1+5)}{\Gamma(1)}$$

$$\text{For } wd \ (j = 2): \frac{\Gamma(5)}{\Gamma(5+73)} \times \frac{\Gamma(1+3)}{\Gamma(1)} \times \frac{\Gamma(1+18)}{\Gamma(1)} \times \frac{\Gamma(1+16)}{\Gamma(1)} \times \frac{\Gamma(1+20)}{\Gamma(1)} \times \frac{\Gamma(1+16)}{\Gamma(1)}$$

$$\text{For } wd \ (j = 3): \frac{\Gamma(5)}{\Gamma(5+50)} \times \frac{\Gamma(1+5)}{\Gamma(1)} \times \frac{\Gamma(1+4)}{\Gamma(1)} \times \frac{\Gamma(1+18)}{\Gamma(1)} \times \frac{\Gamma(1+7)}{\Gamma(1)} \times \frac{\Gamma(1+16)}{\Gamma(1)}$$

$$BDe\left(\frac{G}{D}\right) = wd \ (j = 1) \times wd \ (j = 2) \times wd \ (j = 3) \approx 8.98 \times 10^{-125}$$

The Tabu Search algorithm iteratively explores modifications to the DAG (adding, deleting, or reversing edges) while avoiding cycles. Each candidate structure is evaluated using the BDe score, and the highest-scoring network is selected.

1.3.6 Conditional Probability Tables (CPTs)

With the learned structure, CPTs are computed:

$$P\left(\frac{X = x_k}{Pa(X) = pa_j}\right) = \frac{N_{ijk}}{N_{ij}}$$

Table S3: Conditional Probability Table (CPT) representing $P\left(\frac{bdam}{wd}\right)$.

bdam	$wd \leq 1\text{ m}$	$wd > 1\text{ m} \& wd \leq 2\text{ m}$	$wd > 2\text{ m}$
[0.000, 0.120)	$\frac{30}{65} = 0.46$	$\frac{3}{73} = 0.04$	$\frac{5}{50} = 0.10$
[0.120, 0.300)	$\frac{19}{65} = 0.29$	$\frac{18}{73} = 0.25$	$\frac{4}{50} = 0.08$
[0.300, 0.528)	$\frac{7}{65} = 0.11$	$\frac{16}{73} = 0.22$	$\frac{18}{50} = 0.36$
[0.528, 0.960)	$\frac{4}{65} = 0.06$	$\frac{20}{73} = 0.27$	$\frac{7}{50} = 0.14$
[0.960, 1.000]	$\frac{5}{65} = 0.08$	$\frac{16}{73} = 0.22$	$\frac{16}{50} = 0.32$

1.3.7 Bayesian Inference using the Markov Blanket

With CPTs established, Bayesian inference is conducted using the Markov Blanket of the target node. The posterior is calculated as:

$$P\left(\frac{X_i}{E}\right) = \frac{P\left(\frac{E}{X_i}\right)P(X_i)}{P(E)}$$

Where:

$P\left(\frac{E}{X_i}\right)$ = likelihood of evidence given X_i .

$P(X_i)$ = prior probability of X_i .

$P(E) = \sum_{X_i} P\left(\frac{E}{X_i}\right)P(X_i)$

Example: For evidence $wd > 2\text{ m}$:

1. Compute marginal probabilities $P(bdam_k)$

Table S4: Marginal probabilities of $bdam$ bins.

bdam	count	$P(bdam_k)$
[0.000, 0.120)	30+3+5=38	$\frac{38}{188} = 0.2021$
[0.120, 0.300)	19+18+4=41	$\frac{41}{188} = 0.2181$
[0.300, 0.528)	7+16+18=41	$\frac{41}{188} = 0.2181$
[0.528, 0.960)	4+20+7=31	$\frac{31}{188} = 0.1649$
[0.960, 1.000]	5+16+16=37	$\frac{37}{188} = 0.1968$
Total	188	

2. Compute conditional probabilities $P\left(\frac{wd=3}{bdam_k}\right)$

Table S5: Conditional probabilities

bdam	$wd > 2\text{ m}$ Count	$P\left(\frac{wd > 2\text{ m}}{bdam_k}\right)$
[0.000, 0.120)	5	$\frac{5}{38} = 0.1316$
[0.120, 0.300)	4	$\frac{4}{41} = 0.0976$
[0.300, 0.528)	18	$\frac{18}{41} = 0.4390$
[0.528, 0.960)	7	$\frac{7}{31} = 0.2258$
[0.960, 1.000]	16	$\frac{16}{37} = 0.4324$
Total	50	

3. Compute $P(wd > 2\text{ m})$

$$P(wd > 2\text{ m}) = \frac{50}{188} = 0.26596$$

4. Compute posterior probabilities:

$$P\left(\frac{bdam_k}{wd > 2\text{ m}}\right) = \frac{P\left(\frac{wd > 2\text{ m}}{bdam_k}\right) \times P(bdam_k)}{P(wd > 2\text{ m})}$$

Table S6: Posterior probabilities

bdam bin	$P\left(\frac{wd > 2\text{ m}}{bdam_k}\right)$	$P(bdam_k)$	$P(wd > 2\text{ m})$	$P\left(\frac{bdam_k}{wd > 2\text{ m}}\right)$
[0.000, 0.120)	$\frac{5}{38} = 0.1316$	$\frac{38}{188} = 0.2021$	0.26596	0.10
[0.120, 0.300)	$\frac{4}{41} = 0.0976$	$\frac{41}{188} = 0.2181$	0.26596	0.08
[0.300, 0.528)	$\frac{18}{41} = 0.4390$	$\frac{41}{188} = 0.2181$	0.26596	0.36
[0.528, 0.960)	$\frac{7}{31} = 0.2258$	$\frac{31}{188} = 0.1649$	0.26596	0.14
[0.960, 1.000]	$\frac{16}{37} = 0.4324$	$\frac{37}{188} = 0.1968$	0.26596	0.32

$$P\left(\frac{bdam_k}{wd > 2\text{ m}}\right) = [0.10, 0.08, 0.36, 0.14, 0.32]$$

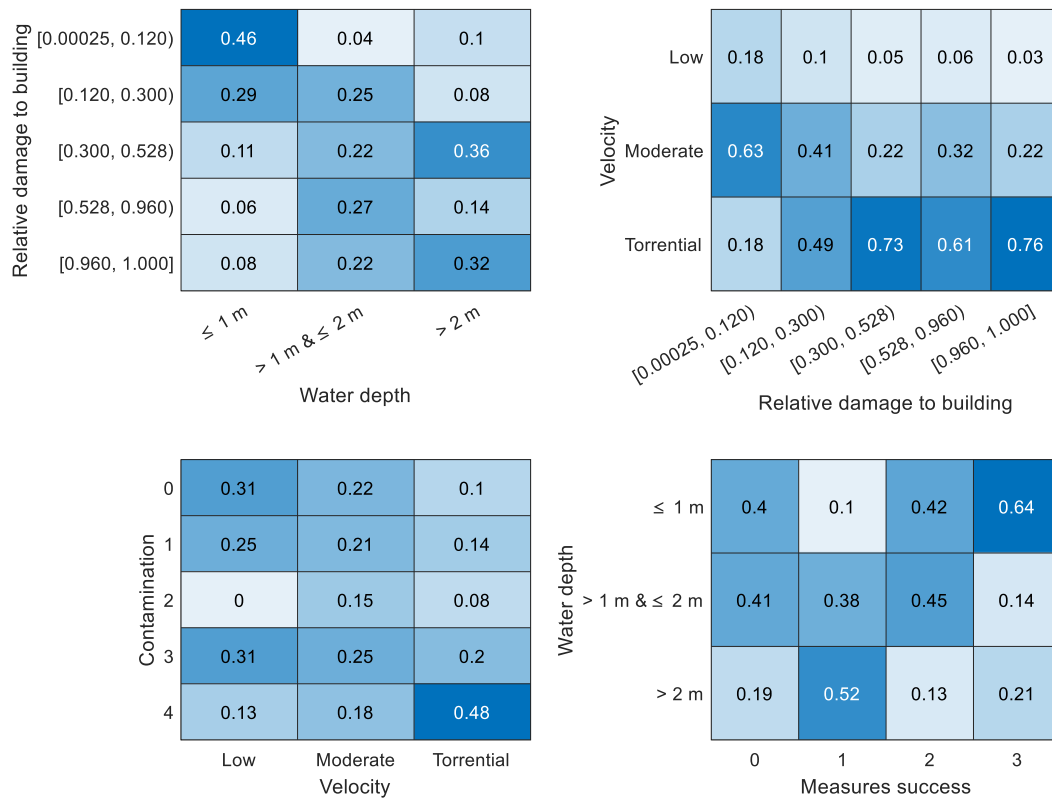


Figure S4: Conditional probability table of companies building Bayesian network (Fig. 7a). X-axis denotes the parent node, while Y-axis denotes the children node within the network.

References:

- Cawley, G. C. and Talbot, N. L. C.: On Over-fitting in Model Selection and Subsequent Selection Bias in Performance Evaluation, *Journal of Machine Learning Research*, 11, 2079–2107, 2010.
- Heckerman, D., Geiger, D., and Chickering, D. M.: Learning Bayesian networks: The combination of knowledge and statistical data, *Mach Learn*, 20, 197–243, <https://doi.org/10.1007/BF00994016>, 1995.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, É.: Scikit-learn: Machine Learning in Python, *Journal of Machine Learning Research*, 12, 2825–2830, 2011.