



# 1 Multi-scale spatial validation and probability calibration of pixel-based 2 landslide susceptibility modeling in the northern Peruvian Andes

3

4 Wendy Quiroz<sup>1</sup>

5

6 <sup>1</sup> Instituto Geofísico del Perú (IGP), Dirección de Ciencias de la Tierra Sólida, Lima, Peru

7 Correspondence to: Wendy Quiroz ([wquiroz@igp.gob.pe](mailto:wquiroz@igp.gob.pe))

8

## 9 Abstract

10 Landslides are recurrent geohazards in Andean regions, causing significant impacts on  
11 infrastructure and local communities. In spatially structured terrains, model reliability hinges on  
12 the definition of pseudo-absence samples and the treatment of spatial dependence during  
13 validation. This study evaluates pixel-based rotational landslide susceptibility in the province of  
14 Huancabamba (Piura, northern Peru) using a Random Forest classifier and seven conditioning  
15 factors derived from a photogrammetric digital elevation model and lithological data at 10 m  
16 resolution.

17 The landslide inventory consists of 25 field-mapped rotational landslides compiled from  
18 geomorphological surveys and high-resolution photogrammetric products. Pseudo-absence  
19 samples were selected outside mapped polygons using a buffered exclusion zone to reduce label  
20 uncertainty, and a balanced sampling scheme (1:1) was adopted. To obtain spatially realistic  
21 performance estimates, model evaluation was conducted using spatial block cross-validation with  
22 block sizes ranging from 600 to 1500 m. This provides a clear view of how spatial partitioning  
23 affects discrimination and calibration, alongside the model's stability throughout the validation  
24 folds.

25 Results show that discrimination performance decreases systematically as spatial block size  
26 increases, indicating that conventional random validation may overestimate predictive capacity  
27 due to spatial autocorrelation. A block size of 900 m provided a compromise between spatial  
28 independence and fold stability. Permutation importance computed under spatially independent  
29 folds identified lithology and elevation as the dominant predictors of rotational landslide  
30 occurrence, followed by aspect and topographic wetness index. Calibration metrics (Brier score  
31 and Expected Calibration Error) indicated moderate but stable reliability of susceptibility scores  
32 across spatial configurations.

33 The resulting susceptibility map shows spatial patterns consistent with the geomorphological  
34 setting and the mapped inventory, with high susceptibility concentrated in steep slopes developed  
35 over weak lithological units. These findings indicate that integrating spatial validation, calibration,  
36 and constrained sampling improves the reliability of pixel-based modelling in this Andean setting.

37 Keywords: landslide susceptibility; Random Forest; spatial cross-validation; pseudo-absence  
38 sampling; Andean

39

## 40 1. Introduction

41 Landslides represent one of the most recurrent and damaging geohazards in mountainous regions  
42 worldwide, producing significant impacts on infrastructure, ecosystems, and population safety. In



43 the Andes, slope instability is strongly controlled by the interaction of steep relief, heterogeneous  
44 lithological conditions, and marked climatic variability, which together favor the recurrence of  
45 mass-movement processes (Hermanns et al., 2012; Guzzetti et al., 2012). In northern Peru,  
46 particularly in the province of Huancabamba (Piura region), rotational landslides frequently affect  
47 transportation corridors, agricultural areas, and expanding peri-urban zones, increasing territorial  
48 vulnerability and socioeconomic losses.

49 Landslide susceptibility mapping (LSM) is widely used in risk reduction strategies and land-use  
50 planning, particularly in mountainous regions (Fell et al., 2008; Van Westen et al., 2008).  
51 Susceptibility assessment aims to estimate the spatial probability of landslide occurrence based on  
52 the relationship between past events and terrain conditioning factors (Corominas et al., 2014). The  
53 shift toward machine learning in landslide susceptibility modeling stems from the need to manage  
54 intricate, nonlinear interactions within high-dimensional multivariate datasets. Among them,  
55 Random Forest has demonstrated robust predictive performance across diverse geomorphological  
56 contexts (Breiman, 2001; Catani et al., 2013; Hong et al., 2019).

57 However, model performance does not depend only on the choice of algorithm. A critical issue in  
58 susceptibility modelling is the evaluation procedure itself. In spatially structured environments,  
59 the assumption of independence between training and validation samples is often violated due to  
60 spatial autocorrelation, which can lead to overly optimistic performance estimates and limited  
61 generalization capability (Chung and Fabbri, 2003; Roberts et al., 2016; Meyer et al., 2019; Valavi  
62 et al., 2019). This issue is particularly relevant for pixel-based susceptibility models, where  
63 neighboring pixels often share similar terrain attributes and may inadvertently introduce spatial  
64 leakage between training and validation subsets.

65 Spatially explicit cross-validation strategies, such as block cross-validation, have been proposed  
66 to mitigate spatial leakage by partitioning data into geographically independent subsets.  
67 Nevertheless, the influence of block size on predictive performance, stability across folds, and  
68 probabilistic reliability remains insufficiently explored in landslide susceptibility studies.  
69 Furthermore, while discrimination metrics such as the area under the ROC curve (AUC) are  
70 commonly reported, considerably less attention has been given to probabilistic calibration and  
71 predictive stability, which are essential for translating susceptibility maps into operational risk  
72 management tools. In susceptibility modelling, predicted values are often interpreted as continuous  
73 susceptibility scores rather than absolute probabilities, particularly when balanced sampling  
74 strategies are adopted. Therefore, in addition to discrimination metrics, calibration measures such  
75 as the Brier score (Brier, 1950; Alvioli et al., 2024) and the expected calibration error (ECE; Guo  
76 et al., 2017) were used to evaluate the reliability of the predicted susceptibility scores.

77 In this context, the study evaluates pixel-based rotational landslide susceptibility in the province  
78 of Huancabamba (northern Peru) under a spatially explicit validation framework designed to  
79 reduce spatial leakage and improve generalization assessment (Roberts et al., 2016; Meyer et al.,  
80 2019). Using Random Forest and terrain conditioning factors derived on a consistent 10 m grid,  
81 we compare spatial block cross-validation schemes to examine their influence on predictive  
82 discrimination and stability across folds. Although spatial cross-validation has been increasingly  
83 recommended in ecological and environmental modelling, systematic multi-scale evaluation of  
84 block size effects on both discrimination and probabilistic calibration remains rare in pixel-based  
85 landslide susceptibility studies.



86 In addition to conventional discrimination metrics, probabilistic calibration metrics are  
87 incorporated to assess reliability, together with an analysis of predictor contributions using  
88 permutation importance computed within spatially independent folds. By focusing on spatial  
89 realism and score reliability, this study aims to improve the methodological consistency of  
90 landslide susceptibility modelling in complex mountainous environments.

## 91 **2. Study area**

92 The study area is located in the province of Huancabamba, Piura region, northern Peru, within the  
93 western flank of the northern Andes (Figure 1). The analyzed sector covers approximately 32 km<sup>2</sup>  
94 and comprises mountainous terrain characterized by steep slopes, deeply incised valleys, and  
95 pronounced relief energy. Elevations range from approximately 1,870 to 2,805 m a.s.l., reflecting  
96 strong altitudinal gradients and complex topographic conditions conducive to slope instability  
97 (Hermanns et al., 2012; Reichenbach et al., 2018).

98 Geologically, the area is composed predominantly of volcanic and sedimentary units of Cenozoic  
99 age, locally overlain by colluvial and alluvial deposits. These lithological formations are frequently  
100 affected by fracturing, weathering, and structural discontinuities, which reduce rock mass strength  
101 and increase susceptibility to gravitational processes. The regional geological framework is  
102 consistent with the national geological map of Peru compiled by INGEMMET (INGEMMET,  
103 2021). Regional assessments report diverse mass movement mechanisms, including rotational  
104 landslides, rockfalls, soil creep, and ground settlements, following the classification of slope  
105 movements proposed by González de Vallejo (2002). The predominance of rotational failures in  
106 the mapped inventory reflects the geomorphological and mechanical characteristics of the local  
107 lithological framework.

108 Climatically, Huancabamba is influenced by a seasonal precipitation regime associated with the  
109 austral summer, with rainfall concentrated between December and March. Intense and prolonged  
110 rainfall events act as major triggering factors for slope instability, particularly in steep terrains  
111 composed of low-cohesion materials (Guzzetti et al., 2008; Gariano and Guzzetti, 2016). In  
112 addition, the region is periodically affected by large-scale climatic anomalies linked to the El  
113 Niño–Southern Oscillation (ENSO), which may substantially increase rainfall intensity and  
114 duration and consequently enhance landslide occurrence (Takahashi et al., 2019; Lavado-Casimiro  
115 et al., 2013).

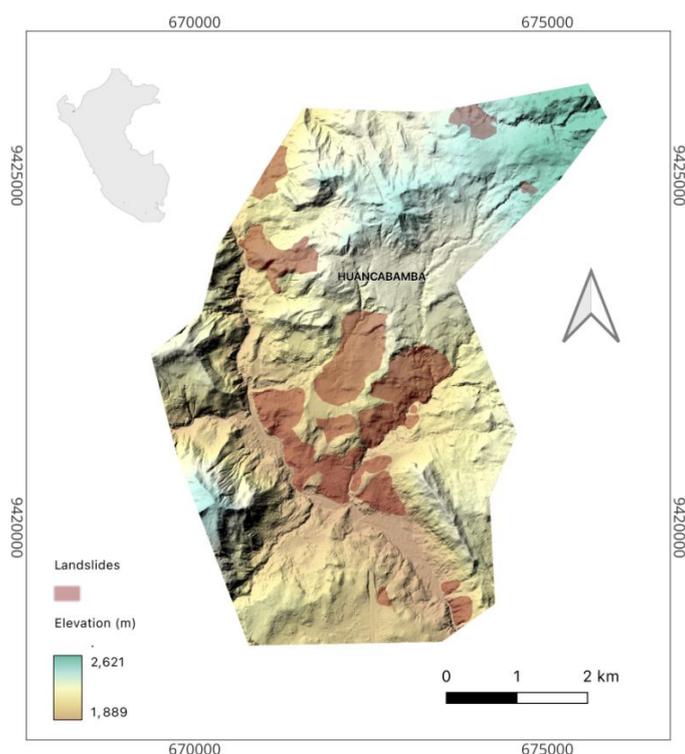
116 Land use within the study area is predominantly agricultural, interspersed with rural settlements  
117 and transportation infrastructure. Anthropogenic modifications, including road cuts, terracing, and  
118 vegetation removal, locally alter drainage patterns and slope geometry, further contributing to  
119 instability processes (Van Westen et al., 2008).

120 Geodynamic investigations conducted by the Instituto Geofísico del Perú (IGP, 2022) documented  
121 active and ancient mass movement processes affecting the urban sector, and GNSS monitoring  
122 between 2018 and 2022 confirmed measurable ground displacements associated with ongoing  
123 slope deformation. The mapped inventory comprises 25 predominantly rotational landslides with  
124 areas ranging from 0.84 to 82.62 ha (mean area  $\approx$  16.7 ha). The characteristic spatial footprint of  
125 the mapped landslides provides a geomorphological scale reference that supports pixel-based  
126 susceptibility modeling at 10 m resolution.



127 In addition, the average event size informed the selection and sensitivity analysis of spatial  
128 validation units, ensuring that block dimensions exceeded the typical landslide footprint while  
129 minimizing spatial leakage.

130 The combination of pronounced topographic gradients, heterogeneous lithological conditions,  
131 seasonal climatic forcing, documented ground deformation, and anthropogenic disturbance makes  
132 the study area an appropriate natural laboratory for evaluating spatially realistic susceptibility  
133 modeling in Andean environments.



134

135 Figure 1. Study area and landslide inventory in the Huancabamba region, northern Peru. The  
136 background shows the digital elevation model (DEM) with hillshade to emphasize the regional  
137 topography. Red polygons represent the landslide inventory used in this study. The inset map  
138 indicates the location of the study area within Peru.

139

### 140 **3. Materials and methods**

141

#### 142 **3.1. Landslide inventory**

143 The landslide inventory used in this study was compiled within the framework of the research  
144 project “Evaluación geofísica y geodinámica de los deslizamientos de tierra que afectan la  
145 seguridad física de la ciudad de Huancabamba, Piura”, conducted by the Instituto Geofísico del  
146 Perú (IGP) between 2018 and 2022. Detailed geomorphological mapping was carried out in 2018  
147 through systematic field surveys supported by high-resolution photogrammetric products. A



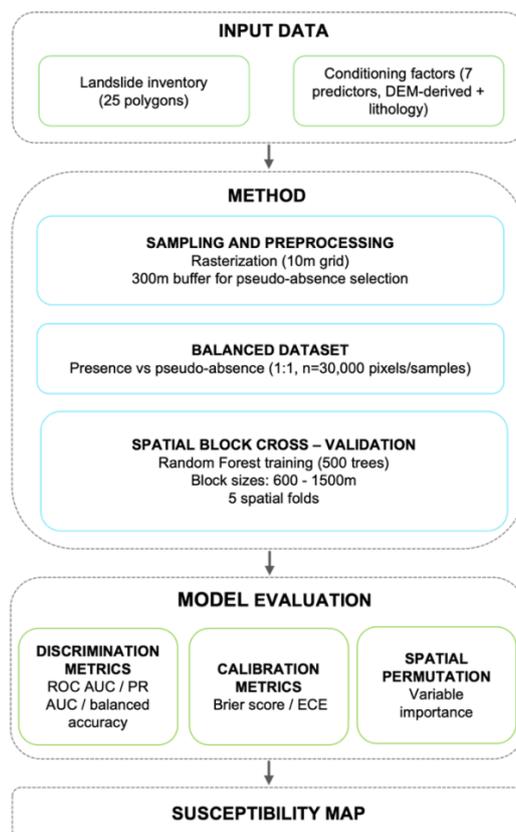
148 consistent and process-based landslide inventory constitutes a fundamental prerequisite for  
149 susceptibility modeling and strongly conditions model validity (Reichenbach et al., 2018).

150 A total of 25 rotational landslide polygons were identified and delineated based on diagnostic  
151 geomorphic features, including main scarps, displaced material, tension cracks, hummocky  
152 topography, and disturbed drainage patterns. These polygons correspond to active instability  
153 sectors documented in regional geodynamic assessments (IGP, 2022). Polygon boundaries were  
154 manually digitized from field-validated photogrammetric products to ensure spatial consistency  
155 and geomorphological reliability, following established best practices in landslide inventory  
156 mapping (Guzzetti et al., 2012).

157 Although the study area exhibits multiple types of mass movement processes—including rotational  
158 landslides, rockfalls, soil creep, and ground settlements (González de Vallejo, 2002)—only  
159 rotational landslides were considered in the present analysis. This decision was adopted to  
160 maintain process homogeneity and reduce class heterogeneity during model training, as  
161 susceptibility modelling is sensitive to the mechanical and geomorphological diversity of slope  
162 processes (Reichenbach et al., 2018). The modelling framework therefore follows a binary  
163 presence–absence scheme, where pixels located within mapped rotational landslide polygons  
164 represent the positive class.

165 Inventory polygons were rasterized on a consistent 10 m grid to match the spatial resolution of the  
166 conditioning factors derived from the resampled DEM and ancillary datasets. All pixels falling  
167 within the rasterized landslide polygons were considered potential presence samples. These were  
168 subsequently integrated into the spatial sampling and cross-validation framework described in  
169 Sect. 3.4.

170 The inventory represents the spatial distribution of mapped active instability at the time of the  
171 2018 survey and provides a geomorphologically consistent reference dataset for pixel-based  
172 susceptibility modelling within the analyzed sector. The overall workflow of the susceptibility  
173 modelling procedure is shown in Figure 2.



174

175 Figure 2. Workflow of the spatially explicit landslide susceptibility modelling framework.  
176 Landslide inventory polygons and terrain conditioning factors were rasterized on a consistent 10  
177 m grid. Pseudo-absence samples were selected outside a 300 m buffer around mapped landslides  
178 to reduce label uncertainty, and a balanced dataset (1:1) was generated. Model evaluation was  
179 conducted using multi-scale spatial block cross-validation (600–1500 m) with five folds to ensure  
180 spatial independence between training and validation data. Discrimination metrics (ROC AUC,  
181 PR AUC, and balanced accuracy), calibration metrics (Brier score and expected calibration error),  
182 and spatial permutation importance were computed within each fold. The final Random Forest  
183 model was trained on the complete dataset to produce the landslide susceptibility map.

184

### 185 3.2. Conditioning factors

186

187 Based on geomorphological understanding of the study area and previous landslide susceptibility  
188 research, seven conditioning factors were selected to characterize terrain morphology and  
189 environmental variability: elevation, slope, plan curvature, slope aspect, topographic wetness  
190 index (TWI), surface roughness, and lithology (Van Westen et al., 2008; Reichenbach et al., 2018).

191

192 Topographic variables were derived from a UAV-based photogrammetric digital elevation model  
193 (DEM) with a native spatial resolution of 8 cm. Although this very high resolution provides  
detailed representation of micro-topographic features, such detail exceeds the characteristic scale



194 of mapped slope failures and may introduce noise unrelated to slope-scale instability processes.  
195 Therefore, the DEM was resampled to a spatial resolution of 10 m using bilinear interpolation  
196 prior to computing terrain derivatives. This resolution ensures compatibility with the  
197 geomorphological footprint of mapped landslides while reducing micro-scale variability.

198 The selected grid size was defined considering the size distribution of the mapped rotational  
199 landslides. The smallest mapped polygon exhibits an equivalent circular diameter of  
200 approximately 103 m, while the mean equivalent diameter is approximately 393 m. At 10 m  
201 resolution, even the smallest event is represented by multiple raster cells, preserving slope-scale  
202 representation while maintaining computational efficiency and methodological consistency with  
203 regional susceptibility studies (Reichenbach et al., 2018).

204 Elevation, slope, plan curvature, and aspect were directly derived from the resampled DEM using  
205 standard terrain analysis tools in QGIS. Slope was computed in degrees and subsequently  
206 converted to radians where required for hydrological indices. The topographic wetness index  
207 (TWI) was calculated using the SAGA GIS implementation based on flow accumulation and slope  
208 derived from the DEM. Surface roughness was computed using the roughness algorithm  
209 implemented in QGIS from the same resampled DEM, which measures local elevation variability  
210 within a moving window. Lithology was derived from the national geological map of Peru  
211 (INGEMMET, 2021), revised and interpreted within the framework of the IGP geodynamic  
212 investigation (IGP, 2022), and rasterized to the 10 m grid using nearest-neighbour resampling in  
213 order to preserve categorical information. The lithological variable captures contrasts in material  
214 composition and mechanical behavior across the study area.

215 All conditioning factors were aligned to a common spatial resolution (10 m), spatial extent, and  
216 coordinate reference system (WGS 84 / UTM zone 17S) prior to sampling and model training,  
217 ensuring consistent pixel-wise predictor extraction and compatibility with the spatial cross-  
218 validation framework.

219

### 220 3.3. Pseudo-absence selection

221

222 In landslide susceptibility modelling, the absence of mapped landslides does not necessarily imply  
223 geomorphological stability, particularly in mountainous environments where latent or unrecorded  
224 instabilities may be present (Reichenbach et al., 2018). Consequently, the definition of pseudo-  
225 absence samples represents a critical methodological step to reduce label variability and prevent  
226 contamination between presence and absence classes (Barbet-Massin et al., 2012; Meyer et al.,  
227 2019; Guo et al., 2024).

228 In this study, pseudo-absence samples were randomly selected from areas located outside mapped  
229 rotational landslide polygons. To reduce the likelihood of including pixels potentially affected by  
230 boundary variability or adjacent slope deformation, a 300 m buffer was applied around all mapped  
231 landslide polygons prior to pseudo-absence extraction. Candidate negative samples were drawn  
232 exclusively from areas beyond this exclusion zone.

233 The buffer distance (300 m) was defined considering the characteristic scale of mapped landslides.  
234 The smallest mapped polygon exhibits an equivalent circular diameter of approximately 103 m,  
235 while the mean equivalent diameter is approximately 393 m. A 300 m exclusion distance therefore  
236 reduces the probability of selecting ambiguous negative samples in the immediate vicinity of



237 mapped instability sectors, while preserving a sufficiently large domain for pseudo-absence  
238 sampling across the study area.

239 To avoid class imbalance effects in the Random Forest classifier and ensure stable discrimination  
240 performance, a balanced sampling scheme was adopted. An equal number of presence and pseudo-  
241 absence samples were used during model training (1:1 ratio), resulting in a total dataset of 30,000  
242 samples. This balanced configuration facilitates interpretation of discrimination metrics and  
243 reduces bias associated with highly skewed class prevalence (Breiman, 2001; Japkowicz and  
244 Stephen, 2002). Because a balanced sampling strategy was adopted, the predicted model outputs  
245 were interpreted as relative susceptibility scores rather than absolute probabilities of landslide  
246 occurrence. This approach is commonly used in susceptibility modelling to improve classifier  
247 stability while maintaining reliable ranking of susceptible areas (Reichenbach et al., 2018).

248 Sampling was conducted prior to spatial block partitioning, and spatial separation between training  
249 and validation data was subsequently enforced through block cross-validation (Sect. 3.4). This  
250 strategy prioritizes spatial separation, label clarity, and methodological transparency, while  
251 allowing the influence of spatial structure on predictive performance to be explicitly evaluated.

252 This procedure ensures that spatial independence between training and validation subsets is  
253 enforced at the block level while preserving a consistent sampling strategy across block  
254 configurations.

### 255 **3.4. Spatial block cross-validation**

256 Conventional random cross-validation assumes independence among samples, an assumption  
257 frequently violated in spatially structured environmental datasets due to spatial autocorrelation  
258 (Roberts et al., 2016; Meyer et al., 2019). In pixel-based landslide susceptibility models,  
259 neighboring pixels often share similar geomorphological and environmental characteristics, which  
260 may lead to spatial leakage between training and validation subsets and consequently to  
261 overoptimistic performance estimates (Samodra et al., 2024).

262 To mitigate this issue, a spatial block cross-validation framework was implemented. The study  
263 area was partitioned into non-overlapping square blocks arranged on a regular grid. Entire blocks  
264 were assigned to either training or validation subsets, ensuring geographical separation between  
265 samples used for model calibration and evaluation. Five spatial folds were generated, and block  
266 assignment to folds was performed using a fixed random seed to ensure full reproducibility while  
267 maintaining spatial contiguity within each block. No clustering or aggregation of neighboring  
268 blocks into larger fold regions was applied; instead, fold membership was determined at the block  
269 level.

270 Block size selection was informed by the characteristic spatial scale of mapped landslides. The  
271 smallest mapped rotational landslide exhibits an equivalent circular diameter of approximately 103  
272 m, while the mean equivalent diameter is approximately 393 m. Block dimensions substantially  
273 larger than these characteristic footprints reduce the probability that individual landslide systems  
274 are spatially split between training and validation subsets.

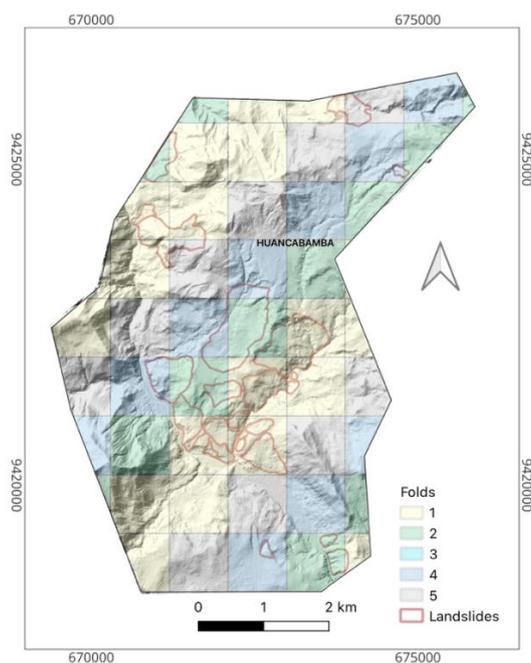
275 Rather than adopting a single block size a priori, a multi-scale sensitivity analysis was conducted.  
276 Several spatial block sizes between 600 and 1500 m were tested in order to evaluate the sensitivity  
277 of model performance to the spatial partitioning scheme.



278 The tested range of block sizes was selected to exceed the mean mapped landslide footprint (~393  
 279 m) while allowing evaluation of model performance across increasing levels of spatial  
 280 independence.

281 For each block configuration, model training was performed exclusively on spatially distinct  
 282 training blocks, and performance metrics were computed on geographically independent validation  
 283 blocks. In configurations where a validation fold contained only one class (presence or absence),  
 284 discrimination metrics requiring two classes (ROC AUC and PR AUC) were not computed for that  
 285 fold and were excluded from aggregated statistics. This conservative treatment prevents artificial  
 286 inflation or distortion of discrimination performance under highly imbalanced spatial partitions.

287 This validation framework allows assessment of discrimination and calibration (Brier score and  
 288 ECE) while maintaining spatial independence during generalization. It provides a spatially  
 289 consistent basis for evaluating pixel-based susceptibility models in the Andes context.



290

291 Figure 3. Spatial block cross-validation configuration using a 900 m block size. The study area  
 292 was partitioned into non-overlapping square blocks assigned to five spatial folds, providing  
 293 geographical separation between training and validation samples.

294

295 Table 1. Spatial cross-validation configuration.

Parameter	Primary configuration	Sensitivity analysis
Block size	900 m	600–1500 m
Number of folds	5	5
Pseudo-absence buffer	300 m	300 m
Class balance	1:1 (balanced)	1:1
Total sample size	30,000	30,000



296

### 297 **3.5. Random Forest model specification**

298

299 Landslide susceptibility was modeled using the Random Forest (RF) algorithm (Breiman, 2001),  
300 an ensemble-based classifier that constructs multiple decision trees through bootstrap aggregation  
301 and random feature selection at each split. RF has been widely adopted in landslide susceptibility  
302 studies due to its ability to capture nonlinear relationships, accommodate multicollinearity among  
303 predictors, and maintain robust predictive performance under complex geomorphological settings  
304 (Catani et al., 2013; Hong et al., 2019; Reichenbach et al., 2018).

305 Random Forest is well-suited for pixel-based susceptibility modelling in mountain terrain, where  
306 predictor variables often display complex spatial patterns and non-linear interactions.

307 The RF model was implemented using 500 trees ( $n_{\text{estimators}} = 500$ ). This value was selected  
308 based on an explicit convergence analysis conducted under the spatial block cross-validation  
309 framework. Model performance metrics (ROC AUC and Brier score) exhibited stabilization  
310 beyond approximately 300–400 trees, while marginal improvements between 500 and larger  
311 ensembles were negligible ( $<0.002$  in ROC AUC). Therefore, 500 trees were retained as a balance  
312 between predictive stability and computational efficiency.

313 Trees were grown without imposing a maximum depth ( $\text{max\_depth} = \text{None}$ ), allowing recursive  
314 partitioning until terminal nodes became pure or contained fewer than the minimum number of  
315 samples required for splitting. The minimum number of samples per leaf was set to one  
316 ( $\text{min\_samples\_leaf} = 1$ ), following common practice in environmental modelling applications  
317 where flexible partitioning of complex predictor space is desirable (Breiman, 2001).

318 Because internal out-of-bag (OOB) error estimates may underestimate generalization error in  
319 spatially structured datasets, predictive performance was evaluated exclusively under spatial block  
320 cross-validation (Sect. 3.4). Block sizes ranging from 600 m to 1500 m were assessed to quantify  
321 the influence of spatial separation on discrimination performance, probabilistic reliability, and  
322 model stability.

323 All models were trained using a fixed random seed to ensure reproducibility of bootstrap sampling  
324 and fold allocation. Given the balanced sampling strategy (1:1 presence to pseudo-absence ratio),  
325 no additional class weighting was applied.

326 Predictor importance was evaluated using two complementary approaches. First, mean decrease  
327 in impurity (MDI) was computed for the full model. Second, and more conservatively, spatial  
328 permutation importance was calculated within each spatial validation fold. This fold-wise  
329 permutation procedure quantifies the reduction in predictive performance when the association  
330 between a predictor and the response is randomly disrupted under spatially independent validation  
331 conditions, thereby providing a more robust assessment of variable contribution in the presence of  
332 spatial autocorrelation.

333 In addition to conventional discrimination metrics (ROC AUC and PR AUC), probabilistic  
334 reliability was explicitly evaluated through the Brier score and the Expected Calibration Error  
335 (ECE). This evaluation allows assessment of both discrimination and calibration metrics,  
336 improving the interpretability of susceptibility scores.



337 Finally, after model evaluation under spatial cross-validation, a final Random Forest model with  
338 the selected hyperparameters was trained using the complete balanced dataset (30,000 samples) to  
339 generate the continuous landslide susceptibility map for the study area.

340

### 341 **3.6. Model evaluation and spatial validation**

342

343 Model performance was evaluated exclusively under the spatial block cross-validation framework  
344 described in Sect. 3.4. For each block configuration, predictive performance metrics were  
345 computed independently for each spatial fold and subsequently summarized using mean and  
346 standard deviation values to quantify inter-fold variability and model stability.

347 Discrimination capacity was assessed using the area under the receiver operating characteristic  
348 curve (ROC AUC) and the area under the precision–recall curve (PR AUC). ROC AUC provides  
349 a threshold-independent measure of separability between presence and pseudo-absence classes,  
350 whereas PR AUC offers complementary information focused on the positive (landslide) class and  
351 is particularly informative under spatially variable class prevalence.

352 To evaluate classification balance at a fixed decision threshold of 0.5 applied to the susceptibility  
353 scores, balanced accuracy was computed. This metric accounts for both sensitivity and specificity  
354 and reduces bias associated with uneven class distribution across spatial folds.

355 Beyond discrimination metrics, calibration of the predicted susceptibility scores was explicitly  
356 evaluated. The Brier score (Brier, 1950) was used to quantify the mean squared difference between  
357 predicted susceptibility scores and observed outcomes, thereby measuring overall probabilistic  
358 accuracy. In addition, the Expected Calibration Error (ECE; Guo et al., 2017, Dormann et al.,  
359 2020) was calculated using 10 prediction bins to evaluate agreement between predicted  
360 probabilities and observed event frequencies across the probability spectrum.

361 Because the modelling dataset was constructed using balanced sampling (1:1 presence and pseudo-  
362 absence), the Random Forest outputs were interpreted as susceptibility scores rather than absolute  
363 event probabilities. Calibration metrics were therefore used to evaluate the reliability of these  
364 scores under spatially independent validation. The results should not be interpreted as  
365 unconditional landslide probabilities at the landscape scale.

366 In spatial folds containing only one class (presence or pseudo-absence), discrimination metrics  
367 requiring two classes (ROC AUC and PR AUC) were not computed and were excluded from  
368 aggregated statistics. This conservative approach prevents artificial inflation or distortion of  
369 performance estimates under highly imbalanced spatial partitions.

370 This evaluation approach assesses discrimination capacity, probabilistic calibration, and predictive  
371 stability under increasing spatial independence.

372

## 373 **4. Results**

374

### 375 **4.1. Multi-scale spatial validation results**

376 Model performance varied systematically as a function of spatial block size. Increasing block size  
377 imposed stricter spatial separation between training and validation samples, thereby reducing  
378 spatial information sharing and providing a more conservative estimate of generalization



379 performance (Figure 4). Block sizes between 600 and 1500 m were evaluated to quantify the effect  
 380 of increasing spatial independence constraints on model performance.

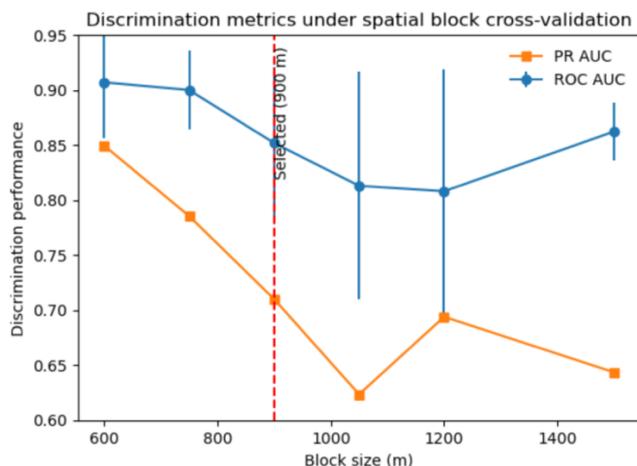
381 Under the 600 m block configuration, the model achieved the highest discrimination performance  
 382 (mean ROC AUC  $\approx 0.91$ ), accompanied by relatively low inter-fold variability. However, this  
 383 configuration likely permits partial spatial dependence between neighboring blocks, potentially  
 384 resulting in optimistic performance estimates.

385 At 750 m block size, discrimination metrics remained high (ROC AUC  $\approx 0.90$ ), while maintaining  
 386 stable fold representation of both classes. Based on the sensitivity analysis, the 900 m block size  
 387 was selected as the primary spatial validation configuration because it provided a balance between  
 388 spatial independence and sufficient sample size within each fold. As block size increased to 900  
 389 m, discrimination performance decreased (ROC AUC  $\approx 0.85$ ), reflecting reduced spatial  
 390 information leakage and more realistic generalization constraints (Figure 4).

391 Further increases in block size (1050–1500 m) led to additional variability across folds. In  
 392 particular, for block sizes  $\geq 1200$  m, some validation folds contained only one class (either presence  
 393 or pseudo-absence), preventing computation of ROC AUC and PR AUC in those folds. This  
 394 behavior reflects the spatial clustering of mapped rotational landslides and highlights the trade-off  
 395 between spatial independence and statistical representativeness.

396 Across block configurations, probabilistic calibration metrics exhibited moderate variation. Brier  
 397 scores ranged approximately between 0.13 and 0.16, while Expected Calibration Error (ECE)  
 398 values remained below 0.18 for all configurations, indicating acceptable reliability of susceptibility  
 399 scores under spatially explicit validation.

400 The 900 m configuration represents a balanced compromise between spatial independence and  
 401 fold stability. It avoids excessive class degeneration observed in larger block sizes while  
 402 substantially reducing potential spatial leakage relative to smaller partitions. For this reason, the  
 403 900 m configuration was retained as the primary spatial validation reference in subsequent  
 404 analyses (vertical dashed line in Figure 4).



405  
 406 Figure 4. Discrimination performance under spatial block cross-validation. Mean ROC AUC and  
 407 PR AUC values are shown across increasing spatial block sizes (600–1500 m). Error bars represent



408 inter-fold standard deviation. The dashed vertical line indicates the 900 m configuration selected  
 409 as the primary spatial validation reference.

410 **4.2. Variable importance under spatial permutation**

411 Variable contribution was evaluated using spatial permutation importance computed within each  
 412 spatial validation fold. Unlike mean decrease in impurity (MDI), which is derived from internal  
 413 tree splits and may be biased in the presence of correlated predictors, permutation importance  
 414 quantifies the reduction in predictive performance when the relationship between a given predictor  
 415 and the response variable is randomly disrupted under spatially independent validation conditions.

416 Under the 900 m block configuration, lithology and elevation consistently exhibited the highest  
 417 importance values across folds. Lithology showed the strongest average contribution (mean  
 418 importance  $\approx 0.23$ ), followed by elevation ( $\approx 0.19$ ), although elevation displayed greater inter-fold  
 419 variability. These results indicate that geological contrasts and large-scale topographic gradients  
 420 exert primary control on the spatial distribution of rotational landslides in the study area.

421 Slope aspect (direction) showed moderate importance ( $\approx 0.016$ ), suggesting that orientation-  
 422 dependent processes—potentially related to differential insolation, moisture retention, or structural  
 423 control—play a secondary but non-negligible role. In contrast, curvature, slope gradient, and  
 424 roughness exhibited low importance values, with roughness occasionally yielding near-zero or  
 425 slightly negative permutation effects, indicating limited independent predictive contribution once  
 426 lithology and elevation are accounted for.

427 Comparable patterns were observed under the 750 m configuration, although importance values  
 428 exhibited greater variability for lithology and elevation. The relative ranking of predictors  
 429 remained stable across block sizes, reinforcing the robustness of the dominant geomorphological  
 430 controls identified.

431 The dominance of lithology highlights the role of material properties and structural discontinuities  
 432 in governing rotational failure development, while elevation likely acts as a proxy for broader  
 433 geomorphic setting and relief energy. The reduced contribution of fine-scale morphometric  
 434 derivatives suggests that, at 10 m resolution and under spatially independent validation, large-scale  
 435 geological and topographic gradients exert stronger predictive influence than local micro-  
 436 topographic variability.

437 Spatial permutation analysis indicates consistent dominance of lithology and elevation across  
 438 independent folds, supporting the geomorphological coherence of the susceptibility model.

439 Table 2. Permutation-based variable importance obtained from the Random Forest model.  
 440 Importance values represent the mean decrease in model performance after random permutation  
 441 of each predictor.

Variable	Mean Importance	Std
Lithology	0.234	0.030
Elevation	0.191	0.062
Aspect	0.016	0.005
Topographic wetness index (TWI)	0.004	0.005
Plan curvature	0.001	0.002
Slope	0.0005	0.003
Surface roughness	-0.0007	0.007



442

443 Importance was computed using permutation importance based on the decrease in model  
444 performance after random shuffling of each predictor.

#### 445 **4.3. Susceptibility map and spatial patterns**

446 The final landslide susceptibility map was generated using the Random Forest model trained on  
447 the complete balanced dataset (30,000 samples) with the selected hyperparameters (500 trees). The  
448 resulting output represents a continuous susceptibility score (0-1), where higher values indicate  
449 higher relative susceptibility (Figure 5). In susceptibility modelling, predicted values are often  
450 interpreted as continuous susceptibility scores rather than absolute probabilities, particularly when  
451 balanced sampling strategies are adopted.

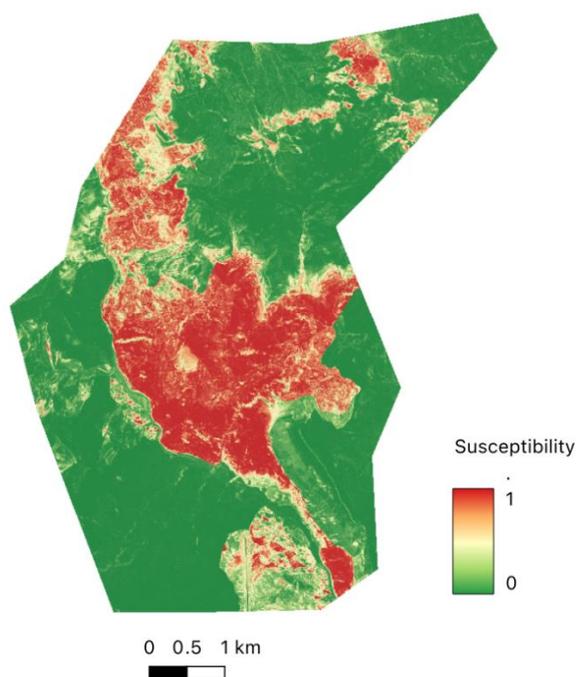
452 Spatially, high-susceptibility zones are predominantly concentrated in sectors characterized by  
453 steep slopes developed over lithological units previously identified as mechanically weaker or  
454 structurally fractured (Figure 5). These areas coincide with mapped rotational landslides and  
455 extend along geomorphological alignments controlled by geological contrasts and relief energy.

456 Intermediate susceptibility values are distributed along transitional terrain zones, often  
457 corresponding to slope breaks, valley flanks, and lithological boundaries. These sectors may  
458 represent areas with favorable conditioning factors but lacking direct evidence of current  
459 instability in the mapped inventory.

460 Low-susceptibility zones are primarily located in relatively gentle slopes, valley bottoms, and  
461 areas underlain by more competent lithological units (Figure 5). The spatial coherence of low-  
462 probability areas supports the internal consistency of the model and suggests that the dominant  
463 conditioning factors exert meaningful geomorphological control.

464 The susceptibility surface exhibits smooth spatial gradients rather than isolated pixel-level  
465 artifacts, reflecting the influence of large-scale predictors such as lithology and elevation. This  
466 pattern is consistent with the spatial permutation importance analysis, which identified geological  
467 and topographic gradients as primary controls of instability.

468 The spatial distribution of predicted susceptibility is consistent with observed geomorphological  
469 patterns and mapped instability sectors.



470

471 Figure 5. Landslide susceptibility map of the Huancabamba study area derived from the Random  
472 Forest model. High susceptibility values (red) coincide predominantly with steep slopes and  
473 weaker lithological units, whereas low susceptibility areas (green) occur in relatively stable terrain.

#### 474 **4.4. Calibration and performance variability**

475 In addition to discrimination metrics, calibration measures were used to evaluate the reliability of  
476 the predicted susceptibility scores. Model stability was assessed through the variability of  
477 performance metrics across spatial cross-validation folds. The Brier score and the expected  
478 calibration error (ECE) were computed to quantify the agreement between predicted scores and  
479 observed landslide occurrence, while the standard deviation of performance metrics across folds  
480 was used as an indicator of model stability under different spatial partitions.

481 Across spatial block configurations, Brier scores ranged between approximately 0.13 and 0.16,  
482 indicating moderate calibration accuracy of the susceptibility scores under spatially independent  
483 validation (Figure 6). Expected Calibration Error (ECE) values remained below 0.18 for all  
484 evaluated block sizes, suggesting acceptable calibration reliability.

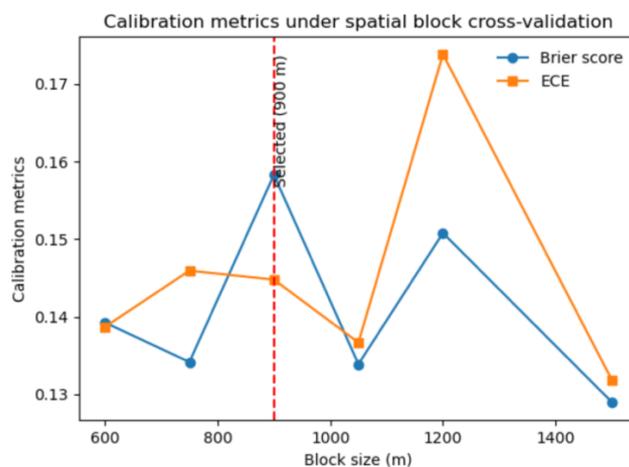
485 Calibration metrics exhibited moderate but non-monotonic variation with increasing spatial block  
486 size (Figure 6). While discrimination performance decreased as spatial independence increased  
487 (Sect. 4.1), calibration indicators remained comparatively stable, without systematic degradation.

488 For block sizes  $\geq 1200$  m, some validation folds contained only a single class, reflecting the spatial  
489 clustering of mapped rotational landslides. Although this structural constraint increased variability  
490 in discrimination metrics, calibration values remained within a relatively narrow range.



491 The observed variability across folds reflects the sensitivity of the model to spatial partitioning  
492 and provides an estimate of the robustness of the susceptibility patterns under different spatial  
493 independence settings.

494



495

496 Figure 6. Calibration performance across spatial block sizes. Brier score and Expected Calibration  
497 Error (ECE) computed under spatial block cross-validation. The dashed line indicates the selected  
498 900 m configuration.

499

## 500 5. Discussion

### 501 5.1 Influence of spatial scale on model generalization

502 Pixel-based susceptibility models are particularly prone to inflated validation scores when spatial  
503 autocorrelation is ignored, because neighbouring pixels share similar predictor values and can  
504 cause spatial leakage under random cross-validation (Roberts et al., 2016; Meyer et al., 2019;  
505 Valavi et al., 2019; Knevels et al., 2023).

506 The multi-scale spatial validation framework implemented in this study reveals a clear trade-off  
507 between discrimination performance and spatial independence. Smaller block sizes (600–750 m)  
508 yielded higher ROC AUC values, suggesting strong apparent predictive capacity. However, such  
509 configurations may allow residual spatial dependence between training and validation subsets,  
510 thereby inflating discrimination metrics through spatial leakage—a well-known issue when  
511 conventional cross-validation assumptions are violated in spatially structured environmental  
512 datasets (Roberts et al., 2016; Meyer et al., 2019; Valavi et al., 2019).

513 As block size increased, discrimination performance systematically decreased and inter-fold  
514 variability increased, consistent with progressively stricter spatial independence constraints and  
515 reduced spatial autocorrelation between training and validation samples (Roberts et al., 2016;  
516 Valavi et al., 2019). For block sizes  $\geq 1200$  m, the occurrence of validation folds containing only  
517 one class further indicates a clustered spatial distribution of rotational landslides within the study  
518 area. This outcome illustrates an important practical limitation of strict spatial partitioning: while  
519 larger blocks increase geographical separation, they may compromise statistical representativeness



520 when the number of independent spatial units becomes small (Valavi et al., 2019; Roberts et al.,  
521 2016).

522 The 900 m configuration represents a geomorphologically informed compromise between spatial  
523 independence and fold stability. This scale exceeds the characteristic footprint of mapped  
524 landslides (mean equivalent diameter  $\approx$  393 m), reducing the likelihood of splitting individual  
525 instability systems across folds while preserving a sufficient number of blocks to maintain class  
526 representation.

527 These results show that model performance depends strongly on the spatial scale used for  
528 validation. Reporting only a single validation configuration, without testing alternative block sizes,  
529 may therefore obscure the influence of spatial autocorrelation on apparent predictive performance  
530 (Meyer et al., 2019; Roberts et al., 2016; Reichenbach et al., 2018).

### 531 **5.2 Geomorphological controls and stability of dominant predictors**

532 Spatial permutation importance analysis consistently identified lithology and elevation as the  
533 dominant predictors of rotational landslide occurrence across spatial block configurations. This  
534 finding highlights the primary role of geological contrasts and large-scale topographic gradients  
535 in controlling slope instability within the study area.

536 Lithology exhibited the strongest and most stable contribution under spatially independent  
537 validation. Unlike internal tree-based importance metrics, spatial permutation importance  
538 evaluates predictor contribution under geographically separated validation subsets, thereby  
539 reducing bias introduced by spatial autocorrelation. The dominance of lithology suggests that  
540 material properties, structural discontinuities, and differential weathering are the most influential  
541 predictors of landslide occurrence in the study area. This result is consistent with previous  
542 landslide susceptibility studies emphasizing the fundamental role of geological framework in slope  
543 instability processes (Van Westen et al., 2008; Reichenbach et al., 2018).

544 Elevation also showed strong predictive contribution, although with greater inter-fold variability.  
545 Rather than acting as a purely geometric descriptor, elevation likely captures broader  
546 geomorphological setting, relief energy, and associated environmental gradients. In Andean  
547 terrain, elevation may integrate multiple process-related controls, including slope steepness  
548 patterns, hydrological gradients, and long-term landscape evolution (Guzzetti et al., 2008; Gariano  
549 and Guzzetti, 2016). Its importance therefore reflects both direct and indirect conditioning  
550 influences.

551 In contrast, local morphometric derivatives such as curvature, slope gradient, and roughness  
552 exhibited limited independent contribution under spatial permutation testing. While such variables  
553 are frequently reported as relevant predictors in conventional susceptibility studies, their reduced  
554 importance under spatially explicit validation suggests that part of their apparent predictive power  
555 may be associated with localized spatial clustering rather than robust generalizable structure. This  
556 observation aligns with recent critiques emphasizing the need to reassess variable importance  
557 under spatially independent frameworks (Meyer et al., 2019; Roberts et al., 2016).

558 The relative ranking of predictors remained stable across block sizes. Even under more rigorous  
559 spatial separation, lithology and elevation continued to dominate. This suggests that the identified  
560 susceptibility patterns reflect consistent geomorphological controls rather than being solely driven  
561 by spatial autocorrelation.

### 562 **5.3 Probabilistic reliability and implications for hazard assessment**



563 While discrimination metrics such as ROC AUC are widely used to evaluate landslide  
564 susceptibility models, they do not assess whether predicted probabilities correspond to observed  
565 event frequencies. In operational risk management contexts, however, the reliability of probability  
566 estimates is as critical as class separability. A model capable of distinguishing presence from  
567 absence may still provide poorly calibrated probability values, potentially misleading decision-  
568 making processes.

569 The evaluation framework adopted in this study incorporated probabilistic reliability metrics,  
570 including the Brier score and the Expected Calibration Error (ECE), under spatially explicit  
571 validation. Across block configurations, calibration metrics indicated moderate and relatively  
572 stable probabilistic agreement, even as discrimination performance decreased under stricter spatial  
573 independence constraints. This behavior suggests that although spatial separation reduces apparent  
574 separability, the model maintains consistent probability structure.

575 Calibration curves revealed minor deviations at higher susceptibility bins, where predicted  
576 probabilities slightly overestimated observed frequencies. Such behavior is not uncommon in  
577 ensemble classifiers trained under balanced sampling schemes and may reflect the interaction  
578 between class balance and spatial clustering. Importantly, the magnitude of calibration deviation  
579 remained limited across spatial scales, indicating that predicted susceptibility values retain  
580 meaningful probabilistic interpretation.

581 From a hazard management perspective, probabilistic reliability enhances the practical value of  
582 susceptibility maps. Well-calibrated susceptibility scores allow thresholds for monitoring or  
583 intervention to be defined based on meaningful risk levels rather than solely on the relative ranking  
584 of susceptibility values. In mountainous environments characterized by spatially clustered  
585 instabilities, ensuring calibration of the predicted scores under spatially independent validation  
586 strengthens confidence in the applicability of susceptibility outputs to territorial planning and  
587 mitigation strategies.

588 Evaluating discrimination, calibration, and spatial stability together provides a more complete  
589 assessment of model estimates.

#### 590 **5.4 Limitations and future research directions**

591 Despite the methodological rigor adopted in this study, several limitations should be  
592 acknowledged.

593 First, the landslide inventory is based on field mapping conducted in 2018 and represents active  
594 instability sectors documented during that survey period. Although geomorphologically  
595 consistent, the inventory does not capture potential future failures or temporal variability in slope  
596 activity. Incorporating multi-temporal inventories could further refine susceptibility calibration  
597 and improve long-term generalization assessment.

598 Second, only rotational landslides were included in the modelling framework in order to preserve  
599 process homogeneity. While this approach reduces class heterogeneity and enhances  
600 interpretability, it limits the applicability of the resulting susceptibility map to other mass  
601 movement types such as rockfalls or shallow translational slides. Future research could explore  
602 multi-class or process-specific modelling frameworks under spatially explicit validation.

603 Third, the balanced sampling strategy (1:1 presence to pseudo-absence ratio) facilitates  
604 discrimination assessment but does not reflect true spatial prevalence of landslides within the study  
605 area. Although probabilistic calibration metrics were evaluated under this configuration, further



606 work could examine the effect of prevalence-aware sampling strategies on probability  
607 interpretation.

608 Finally, the spatial block cross-validation approach, while reducing spatial leakage, is inherently  
609 sensitive to the number and configuration of available spatial units. In relatively small study areas  
610 (~32 km<sup>2</sup>), strict spatial independence constraints may reduce fold representativeness, as observed  
611 in larger block configurations. Expanding the framework to broader regional scales would allow  
612 evaluation of model behavior under increased spatial heterogeneity.

613 Future research incorporating multi-temporal data broader spatial coverage, and additional  
614 physically based predictors could improve evaluation of susceptibility models under spatially  
615 independent validation.

616

## 617 **6. Conclusions**

618 This study evaluated pixel-based rotational landslide susceptibility in the Huancabamba sector of  
619 the northern Peruvian Andes under a spatially explicit validation framework designed to reduce  
620 spatial leakage and provide realistic estimates of model generalization.

621 Model performance was strongly dependent on spatial partition scale. Discrimination decreased  
622 systematically from mean ROC AUC  $\approx 0.91$  at 600 m block size to  $\approx 0.85$  at 900 m, highlighting  
623 the influence of spatial autocorrelation on apparent predictive capacity. Larger block sizes ( $\geq 1200$   
624 m) increased inter-fold variability and led to class imbalance in some folds, illustrating the trade-  
625 off between spatial independence and statistical representativeness in relatively small study areas  
626 (~32 km<sup>2</sup>). A 900 m block size was retained as a compromise, exceeding the mean landslide  
627 footprint (~393 m) while preserving fold stability.

628 Spatial permutation importance computed under geographically independent folds consistently  
629 identified lithology (mean importance  $\approx 0.23$ ) and elevation ( $\approx 0.19$ ) as dominant predictors. Their  
630 stable ranking across block configurations indicates that large-scale geological contrasts and  
631 topographic gradients exert primary control on the spatial distribution of rotational failures in the  
632 study area. In contrast, fine-scale morphometric derivatives showed limited independent  
633 contribution once spatial dependence was constrained.

634 Beyond discrimination metrics, calibration assessment indicated moderate but stable reliability of  
635 susceptibility scores across spatial configurations (Brier score  $\approx 0.13$ – $0.16$ ; ECE  $< 0.18$ ). Although  
636 discrimination declined under stricter spatial separation, calibration remained comparatively  
637 consistent, suggesting that predicted scores preserve meaningful probabilistic structure under  
638 spatially independent validation.

639 These results demonstrate that validation scale substantially affects both apparent discrimination  
640 and stability of pixel-based susceptibility models. Reporting a single validation configuration may  
641 therefore obscure the influence of spatial autocorrelation on performance estimates. Multi-scale  
642 spatial cross-validation provides a transparent means of quantifying this effect and selecting a  
643 validation scale consistent with geomorphological process dimensions.

644 The modelling approach developed in this study is particularly suited to rotational landslides under  
645 similar geomorphological conditions. Extension to other mass movement types would require  
646 process-tailored sampling strategies and predictor selection.

647 Data availability.



648 The datasets used in this study include a landslide inventory and several terrain conditioning  
649 factors derived from remote sensing products and thematic cartography. Because some of the raster  
650 layers are relatively large, the complete dataset is not hosted in the code repository. The data  
651 supporting the results presented in this study can be obtained from the author upon reasonable  
652 request.

653 Code availability.

654 The notebook and scripts used for data preprocessing, Random Forest modelling, spatial block  
655 cross-validation, and performance evaluation are maintained in a GitHub repository prepared for  
656 this study. The repository can be made available upon request during the review process and will  
657 be made publicly accessible after acceptance.

658 Author contributions.

659 Wendy Quiroz designed the study, prepared the datasets, developed the methodology, carried out  
660 the modelling, analysed the results, and wrote the manuscript.

661 Competing interests.

662 The author declares that there are no competing interests.

663 Acknowledgements.

664 The author gratefully acknowledges the support of the Instituto Geofísico del Perú (IGP) during  
665 the development of this research and thanks colleagues and collaborators for helpful discussions  
666 that contributed to the preparation of this work.

667

## 668 **References**

669 Alvioli, M., et al.: A benchmark dataset and workflow for landslide susceptibility zonation, *Earth-*  
670 *Science Reviews*, 258, 104927, <https://doi.org/10.1016/j.earscirev.2024.104927>, 2024.

671 Barbet-Massin, M., Jiguet, F., Albert, C. H., and Thuiller, W.: Selecting pseudo-absences for  
672 species distribution models: how, where and how many?, *Methods Ecol. Evol.*, 3, 327–338,  
673 <https://doi.org/10.1111/j.2041-210X.2011.00172.x>, 2012.

674 Breiman, L.: Random Forests, *Mach. Learn.*, 45, 5–32, <https://doi.org/10.1023/A:1010933404324>,  
675 2001.

676 Brier, G. W.: Verification of forecasts expressed in terms of probability, *Mon. Weather Rev.* 78, 1–  
677 3, 1950

678 Catani, F., Lagomarsino, D., Segoni, S., and Tofani, V.: Landslide susceptibility estimation by  
679 Random Forests technique: sensitivity and scaling issues, *Nat. Hazards Earth Syst. Sci.*, 13, 2815–  
680 2831, <https://doi.org/10.5194/nhess-13-2815-2013>, 2013.

681 Chung, C. J. F. and Fabbri, A. G.: Validation of spatial prediction models for landslide hazard  
682 mapping, *Nat. Hazards*, 30, 451–472, <https://doi.org/10.1023/B:NHAZ.0000007172.62651.2b>,  
683 2003.

684 Corominas, J., van Westen, C., Frattini, P., Cascini, L., Malet, J.-P., Fotopoulou, S., Catani, F., Van  
685 Den Eeckhaut, M., Mavrouli, O., Agliardi, F., Pitilakis, K., Winter, M. G., and Pastor, M.:  
686 Recommendations for the quantitative analysis of landslide risk, *Bull. Eng. Geol. Environ.*, 73,  
687 209–263, <https://doi.org/10.1007/s10064-013-0538-8>, 2014.

688 Dormann, C. F., Calabrese, J. M., Guillerá-Arroita, G., Matechou, E., Bahn, V., Bartoń, K., Beale,  
689 C. M., Ciuti, S., Elith, J., Gerstner, K., Guisan, A., Keil, P., Lahoz-Monfort, J. J., Pollock, L. J.,  
690 Reineking, B., Roberts, D. R., Schröder, B., Thuiller, W., Warton, D. I., Wintle, B. A., Wood, S.



- 691 N., and Hartig, F.: Calibration of probability predictions from machine-learning and statistical  
692 models, *Global Ecology and Biogeography*, 29, 760–776, <https://doi.org/10.1111/geb.13070>, 2020.
- 693 Fell, R., Corominas, J., Bonnard, C., Cascini, L., Leroi, E., and Savage, W. Z.: Guidelines for  
694 landslide susceptibility, hazard and risk zoning for land-use planning, *Eng. Geol.*, 102, 85–98,  
695 <https://doi.org/10.1016/j.enggeo.2008.03.014>, 2008.
- 696 Gariano, S. L., and Guzzetti, F.: Landslides in a changing climate, *Earth-Science Reviews*, 162,  
697 227–252, <https://doi.org/10.1016/j.earscirev.2016.08.011>, 2016.
- 698 Gonzáles de Vallejo, L., Ferrer, M., Ortuño, L., and Otero, C.: *Ingeniería geológica*, Pearson  
699 Educación, Madrid, 2002.
- 700 Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q.: On Calibration of Modern Neural Networks,  
701 in: *Proceedings of the 34th International Conference on Machine Learning (ICML)*, 2017.  
702 Available from <https://proceedings.mlr.press/v70/guo17a.html>.
- 703 Guo, Z., Tian, B., Zhu, Y., et al.: How do the landslide and non-landslide sampling strategies  
704 impact landslide susceptibility assessment? — A catchment-scale case study from China, *Journal*  
705 *of Rock Mechanics and Geotechnical Engineering*, 16, 877–894,  
706 <https://doi.org/10.1016/j.jrmge.2023.07.026>, 2024.
- 707 Guzzetti, F., Peruccacci, S., Rossi, M., and Stark, C. P.: The rainfall intensity–duration control of  
708 shallow landslides and debris flows: an update, *Landslides*, 5, 3–17,  
709 <https://doi.org/10.1007/s10346-007-0112-1>, 2008.
- 710 Guzzetti, F., Mondini, A. C., Cardinali, M., Fiorucci, F., Santangelo, M., and Chang, K.-T.:  
711 Landslide inventory maps: new tools for an old problem, *Earth-Sci. Rev.*, 112, 42–66,  
712 <https://doi.org/10.1016/j.earscirev.2012.02.001>, 2012.
- 713 Hermanns, R. L., Oppikofer, T., Anda, E., Blikra, L. H., Crosta, G. B., and Dahle, H.: Landslide  
714 dams in the Andes of Peru – hazard and risk assessment, *Landslides*, 9, 157–172,  
715 <https://doi.org/10.1007/s10346-011-0300-1>, 2012.
- 716 Hong, H., Pradhan, B., Xu, C., and Tien Bui, D.: Spatial prediction of landslide hazard using  
717 Random Forest, Gradient Boosting Machine, and Naïve Bayes Tree models, *Geomorphology*, 330,  
718 15–29, <https://doi.org/10.1016/j.geomorph.2018.12.021>, 2019.
- 719 INGEMMET: *Mapa geológico del Perú, escala 1:1 000 000*, Instituto Geológico, Minero y  
720 Metalúrgico, Lima, 2021.
- 721 Instituto Geofísico del Perú (IGP): *Huancabamba: Evaluación geofísica y geodinámica de los*  
722 *deslizamientos de tierra que afectan la seguridad física de la ciudad de Huancabamba*, Institutional  
723 report, 2022.
- 724 Japkowicz, N., and Stephen, S.: The class imbalance problem: A systematic study, *Intelligent Data*  
725 *Analysis*, 6(5), 429–449, <https://doi.org/10.3233/IDA-2002-6504>, 2002.
- 726 Knevels, R., Petschko, H., Proske, H., Leopold, P., Mishra, A. N., Maraun, D., and Brenning, A.:  
727 Assessing uncertainties in landslide susceptibility predictions in a changing environment (Styrian  
728 Basin, Austria), *Nat. Hazards Earth Syst. Sci.*, 23, 205–229, [https://doi.org/10.5194/nhess-23-205-](https://doi.org/10.5194/nhess-23-205-2023)  
729 2023, 2023.
- 730 Lavado-Casimiro, W. S., Felipe, O., Silvestre, E., and Bourrel, L.: ENSO impact on hydrology in  
731 Peru, *Advances in Geosciences*, 33, 33–39, <https://doi.org/10.5194/adgeo-33-33-2013>, 2013.
- 732 Meyer, H., Reudenbach, C., Wöllauer, S., and Nauss, T.: Importance of spatial predictor variable  
733 selection in machine learning applications – moving from data reproduction to spatial prediction,  
734 *Ecol. Modell.*, 411, 108815, <https://doi.org/10.1016/j.ecolmodel.2019.108815>, 2019.



- 735 Reichenbach, P., Rossi, M., Malamud, B. D., Mihir, M., and Guzzetti, F.: A review of statistically-  
736 based landslide susceptibility models, *Earth-Sci. Rev.*, 180, 60–91,  
737 <https://doi.org/10.1016/j.earscirev.2018.03.001>, 2018.
- 738 Roberts, D. R., Bahn, V., Ciuti, S., Boyce, M. S., Elith, J., Guillera-Arroita, G., Hauenstein, S.,  
739 Lahoz-Monfort, J. J., Schröder, B., Thuiller, W., Warton, D. I., Wintle, B. A., Hartig, F., and  
740 Dormann, C. F.: Cross-validation strategies for data with temporal, spatial, hierarchical, or  
741 phylogenetic structure, *Ecography*, 40, 913–929, <https://doi.org/10.1111/ecog.02881>, 2016.
- 742 Samodra, G., Ngadisih, and Nugroho, F. S.: *Benchmarking data handling strategies for landslide*  
743 *susceptibility modeling using random forest workflows*, *Artificial Intelligence in Geosciences*, 5,  
744 100093, <https://doi.org/10.1016/j.aiig.2024.100093>, 2024.
- 745 Takahashi, K., and Martínez, A. G.: The very strong coastal El Niño in 1925 in the far-eastern  
746 Pacific, *Climate Dynamics*, 52, 7389–7415, <https://doi.org/10.1007/s00382-017-3702-1>, 2019.
- 747 Valavi, R., Elith, J., Lahoz-Monfort, J. J., and Guillera-Arroita, G. BlockCV: An R package for  
748 generating spatially or environmentally separated folds for k-fold cross-validation of species  
749 distribution models, *Methods Ecol. Evol.*, 10, 225–232, <https://doi.org/10.1111/2041-210X.13107>,  
750 2019.
- 751 Van Westen, C. J., Castellanos, E., and Kuriakose, S. L.: Spatial data for landslide susceptibility,  
752 hazard, and vulnerability assessment: An overview, *Eng. Geol.*, 102, 112–131,  
753 <https://doi.org/10.1016/j.enggeo.2008.03.010>, 2008.