

1 ***Interactive comment on “Assessing the quality of landslide***
2 ***susceptibility maps – case study Lower Austria” by H.***
3 ***Petschko et al.***

4 ***Authors response on Referee #1 comments***

5 **H. Petschko¹, A. Brenning², R. Bell¹, J. Goetz^{1,2}, T. Glade¹**

6 [1] Department of Geography and Regional Research, University of Vienna, Austria

7 [2] Department of Geography and Environmental Management, University of Waterloo,
8 Ontario N2L 3G1, Canada

9

10 Correspondence to: H. Petschko (helene.petschko@univie.ac.at)

11

12 **Reply to the specific comments**

13 We numbered the comments given by the referee, as some comments refer to similar sections
14 which were changed in the revised manuscript. This numbering should assist to make the
15 cross-references within this reply easier to follow.

16 The comments by the referee are in presented in bold, whereas the reply is in normal font
17 type. Please notice that some comments were split up into two comments if the original
18 comment was referring to different lines or aspects. This is indicated with “...” at the end and
19 at the beginning of the comments concerned.

1 *1. p1004 l4ff: In this paragraph the authors motivate their strategy of k-fold cross-*
2 *validation. I feel that, in the introduction chapter, the generic strategy of 1-fold cross-*
3 *validation (or "single hold-out") should be explained more generally, for the less*
4 *experienced reader. In the methods section, the approach can (and must) be explained*
5 *more in detail. Specifically, I suggest to briefly explain the term "single hold-out" as model*
6 *estimation-validation-strategy in one short sentence instead of simply mentioning "single-*
7 *holdout model performance measures".*

8 The suggested brief explanation of the single-hold out performance assessment was included
9 in section 2 of the revised manuscript to make the difference to the *k*-fold cross-validation
10 clearer for the less experienced reader.

11 Added to page 1006, line 7

12 In single hold-out methods the data set is split in one single training and test sample.
13 The training sample is used to fit the model and the test sample is used to determine
14 the model performance. This results in a single estimate of the performance measure
15 (e.g. one single AUROC value) without providing a measure of precision of the
16 estimator. The estimate depends on the (random) sample used for modelling the
17 susceptibility and testing the model's performance, which may itself have "peculiar"
18 random characteristics that would be different for a different test set. Repeated *k*-fold
19 cross-validation solves this problem by using, one after another, different subsets or
20 partitions of the data set as test and training sets, thus effectively using the entire data
21 set for performance estimation (Brenning, 2012a, 2012b). In addition, repeating this
22 procedure for different data partitioning reduces sampling variability and allows for
23 the determination of the precision of the performance estimator (see section 5.3 for
24 details).

1 **2. p1004 19: The abbreviation AUROC is used here without prior explanation what**
2 **AUROC means and, even more importantly, what AUROC IS (the area under a receiver**
3 **operating curve, and as such one possible concept to measure of model quality). As model**
4 **validation is explained in the methods section, I suggest that the authors speak, less**
5 **specifically, of "a range of possible validation outcomes instead of only one single,**
6 **"random" outcome...". Alternatively, the validation concept of ROC and the quality**
7 **measure of AUROC would have to be explained prior to using the abbreviation.**

8 According to your suggestion the term AUROC (area under a receiver operating characteristic
9 curve) is explained earlier in the manuscript and introduced as one possible performance
10 measure. For the detailed explanation and interpretation of the AUROC value we refer in the
11 introduction to the methods section (5.3). Furthermore, the sentence is altered according to
12 your suggestion now being less specific. Generally, this paragraph was moved to section 2 of
13 the revised manuscript. The first occurrence of the abbreviation AUROC is on page 1006 in
14 line 4:

15 Added to page 1006, line 4

16 Among the performance estimation techniques and measures, cross-validation using a
17 single hold-out method and the area under receiver operating characteristic curve
18 (AUROC) value based on ROC plots (e.g. Brenning, 2005; Beguería, 2006; Frattini et
19 al., 2010) are usually applied (i.e. Chung and Fabbri, 1999, 2003; Fabbri et al., 2003;
20 Remondo et al., 2003; Brenning, 2005; Beguería, 2006; Frattini et al., 2010; Rossi et
21 al., 2010).

22 **3. p1010 112f: a comment: The question is also to what extent missing data (from**
23 **blurred, reworked, removed landslides) influences the result of a susceptibility model. If the**
24 **latter is good, the locations of previously existing but now invisible landslide should be**
25 **"predicted" as susceptible. This cannot be checked, but: cross-validation (which reserves**
26 **part of the inventory for evaluation purposes) compared to model goodness-of-fit (i.e. the**
27 **model estimated from and applied to the same area/sample) should give an idea of how an**
28 **incomplete inventory affects model quality.**

29 We are thankful for the correct and helpful comment. This was included in the discussion of
30 the results of the study:

31 Added to page 1026, line 14

1 The implications of an incomplete inventory on the model performance (shown by the
2 AUROC value) were estimated by performing the repeated *k*-fold cross-validation
3 using training and test sample. The results show rather high AUROC values for most
4 modelling domains, which indicates that even with an incomplete inventory (training
5 sample) the prediction of landslides of the test sample was successful for most cases.
6 However, sample size is of importance for the model performance. For the discussion
7 of this please see section 7.4.

8 **4. *p1010 l25ff: Please add some explanation on the possible physical meaning of the***
9 ***DTM-derived variables. For some, it is more obvious (slope) than for others (catchment***
10 ***height - climatic proxy for precipitation ? slope aspect - orientation relative to bedding?)***

11 We understand the lack of description of the DTM derived variables. We included a short
12 description of each variable in the revised manuscript.

13 Added to page 1011, line 6

14 Van Westen et al. (2008) had discussed the relevance of most of these terrain
15 attributes to landslide susceptibility. The relationship of slope angle with landslide
16 activity is well known from general slope stability literature (e.g. Crozier, 1986). Slope
17 aspect can be used as a proxy for bedding orientation. It may also reflect differences in
18 intensity of solar radiation, which controls the local temperature and evaporation and
19 therefore soil moisture (van Westen et al., 2008). Curvature represents convex and
20 concave surfaces related to local morphology (3x3 grid cells). The topographic
21 wetness index was used as a proxy for the soil moisture and ground water level (Beven
22 and Kirkby, 1979; Seibert et al., 2007). The position of the landslide on the slope and
23 the distance from the ridge was represented by the variable catchment height. This was
24 calculated for first order catchments. The catchment area was calculated for the sub-
25 catchments and gives a local representation of the contributing area. Convergence
26 indices were calculated to represent the slope morphology on two different scales by
27 using two different window sizes for the calculation (10m and 50m). Positive values
28 indicate ridges while negative values indicate local depressions.

29 **5. *p1012 l25ff: Yes, tectonic faults are known to influence landslide (in a general***
30 ***sense) activity. In the present study, however, there is a focus on earth (and debris) slides***

1 *(p1008 l17), which means a granulometry of >80 (earth) or <80 (debris) per cent sand and*
2 *finer. I wonder to what degree tectonics influence this type of process compared to, for*
3 *example, lithology, degree of weathering, existence of cover beds, climate etc.*

4 The occurrence of earth and debris slides is very much dependent on the availability of
5 unconsolidated rock or subsequently formed soil of different thickness. The presence of
6 tectonic fault lines is together with the parent material (lithology), climate, topography,
7 vegetation, fauna and time one favouring factor for soil formation (Blume et al., 2010). The
8 closer to a tectonic fault line, the stronger the material is mechanically stressed and
9 fragmented. Therefore, the soil formation process might be faster compared to solid rock.
10 This can also be observed in the field as stated by the Geological Survey of Lower Austria
11 (pers. Comm. Schweigl, 2013). Although earthquakes occur in Lower Austria the influence of
12 earthquakes on landslides was not analysed yet (Hammerl and Lenhardt, 1997). However,
13 earthquakes along active tectonic lines are considered as a possible natural trigger of
14 landslides (probably rather rock slides or rock falls than slides; Schwenk, 1992).

15 The variable “Euclidian distance to tectonic fault line” was selected in 9 out of 16 models
16 using the entire landslide data for fitting the model. This shows that the data has some
17 explanatory power for the occurrence of earth and debris slides. However, we found that the
18 contribution of this variable to the overall model was marginal compared to other variables
19 such as slope angle. No changes made.

20 **6. *p1014 l13ff: "landslide points" are mentioned here and in some other paragraphs,***
21 ***while throughout most of the papers, landslide cells are addressed. Please homogenise***
22 ***terminology. Generally speaking, "points" and "cells" are spatial units to which the***
23 ***sampling and the modelling are applied, and in "reality", you use raster cells, not point***
24 ***objects.***

25 The authors acknowledge the main used terms and changed the wording to landslide cells in
26 all cases throughout the revised manuscript.

1 **7. p1014 l24: Pls define sampling rate (as the number of sampled cells per unit area)**
2 **as opposed to sample size. ...**

3 The reviewer correctly points out a source of confusion in the manuscript. The authors added
4 some sentences for clarification of the difference of sampling rate to sample size in the
5 revised manuscript:

6 Added to page 1014, line 21

7 The sampling of non-landslide cells for the entire study area was based on a density of
8 2% of all grid cells. An equal number of cases and controls (1:1) was used for each
9 model fitted; the landslide samples were matched to an equal number of randomly
10 selected non-landslide samples. This gives the sample size in the respective modelling
11 domain.

12 It was necessary to adjust each model's raw predictions based on the corresponding
13 sampling rate to account for the general relative landslide susceptibility of each
14 modelling domain. We adjusted the prediction by considering the sampling rate (τ_0/τ_1)
15 of each lithology unit, using Eq. 3,

$$16 \text{ odds}^*(x) = t_0/t_1 \cdot \text{odds}(x) \quad (3)$$

17 where,

$$18 t_0 = \text{number of sampled non-landslide cells}/\text{total number of non-landslide cells} \quad (4)$$

19 and

$$20 t_1 = \text{number of sampled landslide cells}/\text{total number of non-landslide cells} \quad (5)$$

21 and $\text{odds}(x)$ is the unadjusted prediction, in our case, based on training a model with a
22 1:1 sampling ratio of landslides to non-landslides.

23 **8. ...Could you also comment on the justification of the 1:1 ratio of landslide to non-**
24 **landslide cells? I suppose it has to do with the binary target variable and the cut-off (of 0.5)**
25 **to distinguish predicted events from non-events.**

26 The decision on a number of cases and controls is very important for the study design.
27 Heckmann et al. (2013) gave a detailed review on which ratios have been chosen in landslide
28 susceptibility modelling and on which grounds (including rare events logistic regression of
29 King and Zeng (2001)). According to these, ratios between 1:2 and 1:5 are recommended to

1 better reflect the full range of values of the explanatory variables. Additionally, it was found
2 that an increase in the ratio of controls to cases larger than four only marginally leads to an
3 increase in precision or statistical power (Ury, 1975; Breslow and Day, 1980). It is rather
4 recommended to increase the number of cases (by widening the study area geographically or
5 temporally) to achieve a better precision (Wacholder et al., 1992). However, the cost of
6 gathering more data on cases (landslides) is rather high in this study, as new landslides would
7 have to be mapped from aerial photographs. As the number of cases and controls is large in
8 most modelling domains the common selection ratio of controls to cases of 1:1 was selected
9 (Breslow and Day, 1980). No changes made.

10 **9. *p1014 l26ff: Why is it necessary to adjust the predictions? The model predicts, in***
11 ***each domain, the probability [0,1] of landslide occurrence, no matter what the sample size***
12 ***or sample rates are. Could you explain this further ? ...***

13 The necessity to adjust the predictions according to the sampling rates of each modelling
14 domain arises due to the need of providing a comparable landslide susceptibility map for the
15 entire study area. Without adjusting the values every modelling domain would show
16 probability values ranging from 0 to 1. However, the general differences in the susceptibility
17 to landslides of each domain as expressed by the sampling rate or landslide density would not
18 be taken into account. Therefore, the probabilities of each modelling domain have to be
19 adjusted according to their sampling rates.

20 For the changes made in the text please refer to our reply on comment 7 and the inserted text
21 in the revised manuscript.

22 **10. *... Furthermore, I could not comprehend the definitions of tau0 and tau1. tau0 is***
23 ***introduced as a "sampling rate for non-landslide points" and defined as the ratio of***
24 ***(land)slide to non-(land)slide points - but this ratio was described earlier as being unity. I***
25 ***understand "sampling rate" as the ratio of sampled cells and the total number of cells in***
26 ***the domain... Your "sampling rate tau1" is defined "for landslide points" and set to 1. This***
27 ***is confusing, because this seems to address the ratio of slide to non-slide cells that was***
28 ***explained earlier; your equation (4) simply becomes tau0*odds(x) because tau1 equals 1...***
29 ***As I understood it, the sampling rates of landslide and non-landslide "points" should be***
30 ***the same in each domain, because of the 1:1 ratio of sampled landslide and non-landslide***
31 ***points. Perhaps tau0 should be the sampling rate (ratio of slide and non-slide pixels to all***

1 *pixels) in the study area, and tau1 this ratio in the domain ? Or is tau0 the ratio of slide to*
2 *non-slide pixels in the study area, and tau1 is the ratio of slide to non-slide points in the*
3 *sample (so tau1 = 1:1 = 1)? To me, an "adjustment" only makes sense if a property of a*
4 *domain (e.g. the ratio of sample size and total size, the ratio of landslide pixels to the total*
5 *area, or the number of landslide pixels) is normalised with the same property of the*
6 *complete study area. Please explain and clarify.*

7 The referee points out an important source of misunderstanding. We mainly identify it as an
8 misunderstanding of the terms sample size and sampling rate and some confusion we created
9 in the equations for τ_0 and τ_1 . The difference between sample size and sampling rate was
10 clarified in our reply to comment 7. The equations of calculating τ_0 and τ_1 were updated
11 and added in the revised manuscript (please refer to comment 7 and the added text in the
12 revised manuscript). In the original manuscript we stated that τ_0 is calculated from the
13 number of slide cells divided by the number of non-slide points. However, the number of
14 sampled non-slide cells has to be divided by the total number of non-slide points. As the
15 sampled number of non-slide cells is identical with the number of slide cells we did not
16 differentiate here between the number of slide cells and sampled number of non-slide cells in
17 the original manuscript. However, we acknowledge that this was rather confusing and
18 therefore we changed this accordingly. Furthermore, we agree that the equation for the
19 calculation of the adjusted odds could be cancelled to $odds^*(x) = \tau_0 \times odds(x)$. However, as
20 future studies might decide to use a subsample of the slide cells we decided to present the
21 entire equation.

22 For the added equations please refer to our reply to comment number 7 and the added text
23 there.

24 *11. p1016 ll6: It is easily understood that the random (non-spatial) partition gives*
25 *different results for every replication, and the strategy of taking the median and IQR of n*
26 *partitions is feasible. Two questions arise for me: (1) did you check if the median and IQR*
27 *of the performance measure are already reasonably stable with 20 replications, or has the*
28 *number of 20 been chosen arbitrarily ? ...*

29 The number of 20 replications is a compromise between a desirable high precision of the
30 cross-validation estimator and acceptable computational complexity given the large sample
31 size and the computational complexity of generalized additive models. Some basic and

1 approximate back-of-envelope calculations based on Figure 4 (in the original manuscript) can
2 give a general idea of the precision of cross-validation estimators. For simplicity we will have
3 a look at the standard error of estimators of mean AUROC (as opposed to the median) across
4 all 20 repetitions. We obtain the following statistical characteristics for domains 3786 (“worst
5 case”: high amount of variation between repetitions) and domain 35 (seems to be more
6 “average”) (IQR=observed sample (i.e. between-repetition) interquartile range; SD = sample
7 standard deviation; SE = standard error, i.e. precision of the estimator of mean AUROC):

8 Domain 3786: IQR approx. 0.20; std.dev. approx. $IQR/1.35 = 0.15$; $SE = SD/\sqrt{100} =$
9 0.033

10 Domain 35: IQR approx. 0.1, std.dev. approx. 0.074, SE approx. 0.017

11 Given the great amount of variation in AUROC among domains (median values from 0.52 to
12 0.98) this precision appears to be highly acceptable, and differences among domains cannot
13 be explained by insufficient estimator precision. Figure 5 (in the original manuscript) also
14 shows very little random variation, which supports this interpretation. We hope that the
15 reviewer will be satisfied with this simplified consideration. No change made.

16 ***12. ... (2) could you add a sentence explaining why/how the k-nearest-neighbour***
17 ***clustering of the "point" coordinates in the spatial partition approach leads to different***
18 ***results for every replication ? Are the k group centroids chosen randomly? Are the***
19 ***resulting spatial units similar in size or is it possible that models are estimated and***
20 ***validated in two partitions of very different size?***

21 The algorithm used here is in fact the *k*-means algorithm, not the k-nearest-neighbour
22 algorithm – apologies for the mistake, which has been corrected. The *k*-means algorithm
23 essentially creates *k* Thiessen polygons. Since the algorithm starts with a random seed of
24 initial cluster centres, different cluster locations and shapes result in each repetition. The
25 resulting spatial units are of similar (comparable) size. Cross-validation ensures that the
26 (approximate) proportion $(k-1)/k$ of the data is used for training, and effectively (over all
27 cross-validation folds and repetitions) all data are used for performance estimation, as
28 indicated in the manuscript.

29 Added to page 1006, line 7

1 Repeated k -fold cross-validation solves this problem by using, one after another,
2 different subsets or partitions of the data set as test and training sets, thus effectively
3 using the entire data set for performance estimation (Brenning, 2012a, 2012b).

4 **13. p1017 l1ff: Why is the IQR an estimated one? It is an empirical measure derived**
5 **from 20 replications of a cross-validation procedure, i.e. of 100 empirical AUROC values**
6 **(see p1016 l18f). ...**

7 Empirical “calculations” are formally referred to as “estimation” in statistics.

8 **14. ... Furthermore: Why does that measure have to be adjusted in order to be a**
9 **measure of transferability? Perhaps you need to explain this more thoroughly, in a more**
10 **step-by-step fashion. ...**

11 The IQR value is influenced by the sample size used for performing the k -fold cross-
12 validation. As the sample size of each modelling domain is very different (104 to 12562 cells)
13 the sampling variability of the AUROC varies. To provide a comparable transferability index
14 for the entire study area this influence of the sampling variability has to be taken into account.
15 Therefore, the IQR was adjusted for the contribution of the standard error SE of the AUROC
16 estimator. Some more explanatory text was added in the revised version of the manuscript.

17 Added to page 1017, line 1

18 The non-spatial and spatial transferability were expressed by the interpretation of the
19 estimated interquartile range (IQR) of the AUROC values resulting from the non-
20 spatial and spatial cross-validation of each modelling domain. The lower the estimated
21 IQR the better we considered the non-spatial and spatial transferability of the model
22 within the modelling domain. Sample size differences among modelling domains
23 result in differences in sampling variability of AUROC estimators, which then has an
24 influence on the IQR of AUROC among cross-validation repetitions. In order to
25 account for this contribution to sampling variability and be able to provide a
26 transferability measure that was comparable among modelling domains, the IQR has
27 to be adjusted according to the sample size.

1 **15. ... Moreover, "Eq (1)" refers to the corresponding equation in the paper of Hanley**
2 **and McNeil, not to Eq(1) in your study... I feel that this equation should be given here, or**
3 **that "Eq(1)" should be removed. ...**

4 According to your suggestions the equation (1) of Hanley and McNeil (1982) was included in
5 the revised manuscript.

6 Added to page 1017, line 5

7 For this purpose, we calculated the approximate standard error SE of the AUROC
8 (AUC) estimator on a test set of N landslide and N non-landslide samples using the
9 equation presented by Hanley and McNeil (1982):

$$10 \quad SE = (AUC \times (1 - AUC) + (N - 1) \times (Q_1 - AUC^2) + (N - 1) \times (Q_2 - AUC^2)) / N^2 \quad (5)$$

11 The quantities Q_1 and Q_2 were calculated from the AUROC (AUC) value as shown
12 by the following equations:

$$13 \quad Q_1 = AUC / (2 - AUC) \quad (6)$$

$$14 \quad Q_2 = 2 \times AUC^2 / (1 + AUC) \quad (7)$$

15 **16. ... Concerning equation (5) for the calculation of the T index: Not knowing "Eq 1**
16 **of Hanley and McNeil", I suspect that the SE of the AUROC is estimated from the standard**
17 **deviation of AUROCs, and will decrease with larger n. I do not see why T should be a better**
18 **indicator of transferability than the empirical IQR, for example. A larger IQR means that a**
19 **model could be very good, but also really bad, while a smaller IQR indicates that the models**
20 **predict similarly well (or poorly). I feel that T should be better justified and explained. I**
21 **understand from the paragraph that you "correct" a non-parametric empirical measure of**
22 **AUROC variability (IQR) by subtracting a parametric, estimated measure of variability**
23 **(that is multiplied by 1.35 to supposedly have the same value as IQR under the assumption**
24 **of normality...). But why ? It may be correct, but it needs more explanation.**

25 We agree on the need for a clarification about the motivation of calculating a transferability
26 index instead of using the estimated IQR values directly as a transferability measure. Please
27 refer to our reply to the comments 14 and 15 as we clarified the questions of the current
28 comment in our reply to the previous comments and in the corresponding changes to the
29 manuscript.

1 **17. p1017 l17ff: assessing the thematic consistency with an index is a good idea.**

2 The authors are thankful for the positive feedback on the thematic consistency index.

3 **18. p1019 l24: The classified susceptibility map is mentioned here, but the classification**
4 **rules are introduced in the results chapter (p1020 l7ff). This should be done before/at the**
5 **first instance when the classified map is introduced (here, or somewhere else in the methods**
6 **section).**

7 We agree with the referee that the methodology of classifying the map has to be introduced in
8 the methods section. However, we would like to refer to page 1015 of the original manuscript
9 where the classification rules of the final susceptibility map were presented in the paragraph
10 starting from line 12. There we point out that we selected the classes according to the
11 percentage of slides contained in each class. The final thresholds of 5% of slides in the lowest
12 susceptibility class and 70% of slides in the highest susceptibility class was a result of testing
13 different thresholds which were checked in the field according to the best geomorphic and
14 planning plausibility.

15 Please refer to page 1015, line 12ff of the original manuscript for more details on the
16 classification rules applied within this study. No changes made.

17 **19. p1021 l10f: Considering that AUROC for the spatial (more meaningful) partition is**
18 **0.53 (very close to useless), the contrast to AUROC=0.79 (acceptable) for random partition**
19 **is very good evidence for the consequences of using over-optimistic validation strategies!**

20 The referee correctly identifies one principal finding of our study. We extended our statement
21 on the performance of the spatial cross-validation in the discussion section of the revised
22 manuscript.

23 Added to page 1028, line 6

24 The median AUROC values estimated by spatial and non-spatial cross-validation were
25 similarly high in this study. However, the median AUROC values and the
26 transferability index clearly showed that non-spatial cross-validation provided a more
27 optimistic or maybe also over-optimistic assessment of the model performance and
28 transferability in contrast to spatial cross-validation.

1 **20. p1021 l13: You mention a (1st quantile) AUROC value of 0.35 - but the AUROC**
2 **only takes values between 0.5 and 1 (see also p1016 l23)!**

3 The authors are thankful for pointing out the source of confusion we created in the original
4 manuscript. We rephrased the sentence on page 1016 in order to be more specific about which
5 values the AUROC can take and which are describing the model's ability to discriminate
6 landslide and non-landslide cells. Values below 0.35 are worse than what would (on average)
7 be expected to be achieved by chance alone, but of course actual estimated values can be
8 worse than that, adding to our argument that resampling-based performance estimation by
9 cross-validation is needed to achieve AUROC estimators of high precision (see comment 11).

10 Added to page 1016, line 22

11 The AUROC takes values between 0 and 1 where a value of 0.5 would be achieved by
12 pure chance agreement between predictions and observations and a value of 1
13 represent perfect discrimination (Brenning, 2005; Guzzetti et al., 2006); however, this
14 may also indicate overfitting. Thus, the AUROC measures the model's ability to
15 discriminate landslide and non-landslide cells (Hosmer and Lemeshow, 2000).

16 **21. p1022 l12ff: "thus the transferability" - does "thus" also apply for n between 200**
17 **and 400 as in line 10f ?...**

18 We agree on the need for more specific wording. The "thus" refers to the increases of the
19 interquartile range at a sample size of about 400 and 200. As the IQR increases the
20 transferability of the model decreases substantially, a trend which becomes stronger the
21 smaller the sample size is.

22 Added to page 1022, line 12

23 ; thus the transferability of the model decreased substantially for sample sizes smaller
24 than 400 (spatial cross-validation) and 200 (non-spatial cross-validation).

1 **22. ... Moreover, can you recommend from your findings a minimum sample size? If so,**
2 **is it related to i) absolute sample size, or ii) to the corresponding sampling rate ? This**
3 **question can possibly be answered with 16 large domains with different landslide**
4 **densities...**

5 As stated in the discussion section (pages 1028 & 1029) both the sample size and the
6 sampling rate have an influence on the transferability and on the thematic consistency of the
7 model as analysed within the k -fold cross-validation. We found that a small sample size and a
8 small sampling rate referring to a large modelling domain with only few landslide
9 occurrences leads to a low transferability and consistency index. Furthermore, the minimum
10 prediction' standard error is lower with a large sample size.

11 However, there are exceptions from this general trend (Tab. 3). Therefore we have to state
12 that the minimum sample size is also always related to the topographic and geotechnical
13 characteristics of the study area. For a domain of similar size and homogeneity in local terrain
14 conditions, a minimum sample size around 400 sample cells is recommended. Therewith the
15 transferability and thematic consistency showed a distinct change using spatial cross-
16 validation at this value (Fig. 5 in the original manuscript). However, the influence of a smaller
17 sample size on the geomorphic plausibility (Bell, 2007) of the susceptibility map is unclear. It
18 might be possible that with a smaller sample size the geomorphic plausibility of the map is
19 lower. Nevertheless, the analysis of this was beyond the scope of this study.

20 We inserted a new small section (7.4) on this in the discussion:

21 Added to page 1030, line 29

22 **7.4 Considerations on sample size**

23 Summarizing the previously discussed findings some considerations on a minimum
24 sample size might be possible. While the transferability index is less strongly related to
25 sample size or sampling rate the thematic consistency index shows a stronger
26 relationship to them. Generally, larger sample sizes and sampling rates result in better
27 thematic consistency and transferability of the model. Furthermore, the minimum
28 standard error of the prediction was lower with larger sample size (Table 3).

29 The effect of a reduced sample size on the median and interquartile range AUROC
30 values was assessed in the Flysch domain. We found that the median AUROC remained
31 satisfactory high but decreased as sample size decreased, while the interquartile range of

1 the AUROC increased. Even with the smallest sample size the model still achieved a
2 good discrimination between landslide and non-landslide cells according to the median
3 AUROC value. Summarizing the results, a minimum sample size with a sum of around
4 400 slide and non-slide cells might be recommended for the methods applied in this
5 study. This size leads to an acceptable transferability and thematic consistency of the
6 model in spatial cross-validation. However, examples from successfully fitting a
7 susceptibility model with smaller sample sizes (10 landslides with 15 cells each in an
8 area of 177km²; Demoulin and Chung, 2007) give a very contrasting result.
9 Furthermore, the sample needs to be substantially complete which might be difficult to
10 estimate for small samples (Malamud et al., 2004). However, increasing the sample size
11 can only be done by enlarging the landslide inventory (e.g. by selecting a larger study
12 area). This is challenging, as in some regions no additional data on landslides (or
13 resources for mapping landslides) might be available.

14 This study showed that the general trends found for sample size and sampling rate do
15 not apply for all modelling domains. Therefore, we highlight that the resulting quality
16 estimates (transferability index, consistency index and prediction uncertainty) might
17 additionally be dependent on a combination of the domain size and the landslide density
18 (landslides per km²). Also, dependencies on local terrain conditions and their
19 homogeneity in the modelling domain might exist.

20 Moreover, the geomorphic plausibility of the susceptibility map has to be analysed.
21 Previous studies found that high performance measures do not always guarantee high
22 geomorphic plausibility of the map (Bell, 2007; Trigila et al., 2013). It might be
23 possible that with a smaller sample size the geomorphic plausibility of the map is lower.
24 However, the influence of a small sample size on the geomorphic plausibility of the
25 susceptibility map is unclear. Nevertheless, analysing this was beyond the scope of this
26 study.

27 **23. p1022 l20f: "one specific random sample and variable selection repetition" - you**
28 **should delete "repetition" (because the results are based on one sample and subsequent**
29 **variable selection).**

30 "Repetition" was deleted in the revised manuscript according to your suggestion.

1 **24. p1022 l27: the possible physical meaning of the variable "catchment height" is not**
2 **explained, neither here nor in the section introducing the variables. See also my comment**
3 **on p1010 l25ff)**

4 Please refer to our authors reply on the referee's comment number 4.

5 **25. p1023 l8: why "on average" ? You have x model runs, and p percent of them**
6 **included the variable.**

7 Changed to "on average over all modelling domains."

8 **26. p1024 l20ff: the propagation of uncertainty to susceptibility classes is a very good**
9 **idea in order not to over-interpret uncertainty while at the same time giving end-users the**
10 **chance of having a closer look where uncertainty crosses the boundary of one or even two**
11 **classes.**

12 The authors are thankful for the positive feedback on our proposed methodology. No changes
13 made.

14 **27. p1028 l23ff: what does "adverse effects" mean? Does that mean that the**
15 **performance measure is biased? or wrong? or that the performance could be better than**
16 **estimated?**

17 This comment refers to criticism of Guzzetti et al. (2006) on the spatial partitioning of the
18 landslide data into training and test sample. Their discussion states "splitting the study area
19 into two adjacent sub-areas can be problematic" as this approach assumes similar
20 characteristics of the explanatory variables (Guzzetti et al., 2006). However, the statement of
21 Guzzetti et al. (2006) is more understood as a word of caution, pointing out the possible
22 pitfalls of spatial partitioning of the study area. Splitting the study area into modelling
23 domains is a major step in the modelling towards reducing possible differences in the
24 characteristics of the explanatory variables. However, as we stated in the discussion section
25 (7.2), it might still be possible that in some samples one terrain condition (e.g. flat areas) is
26 overrepresented relative to others.

27 Added to page 1028, line 25

28 Here, similar characteristics of the explanatory variables in both training and test
29 sample are assumed and necessary (Guzzetti et al., 2006). If this assumption is not met

1 by the data (e.g. a rock type or land use class is missing in the test sample) the transfer
2 of the fitted model to the test sample and the estimation of the model performance are
3 difficult (or impossible) (Guzzetti et al., 2006). In our study some model domains
4 might have high contrast between stable (e.g., large flat areas) and unstable (e.g., steep
5 areas) terrain which gives potential for greater variation of sampled terrain conditions;
6 it may be possible that in some samples one terrain condition is overrepresented
7 relative to others.

8 **28. p1028 l28: I slightly doubt that serious over- or underrepresentation is possible with**
9 **large samples in the order of hundreds to tens of thousands of pixels. Perhaps in very**
10 **inhomogeneous study areas - but that is being dealt with in your approach by establishing**
11 **domains (at least with respect to lithology).**

12 The total sample size (double the number of landslides) is rather small compared to the total
13 number of cells available in each modelling domain (represented by the area in km² of the
14 respective modelling domain in Table 1). Therefore, we did not want to exclude the
15 possibility of over- or underrepresentation of terrain conditions in one sample. However, the
16 detailed analysis of this was beyond the scope of this study. No changes made.

17 **29. p1029 l27f: Does an underestimation not occur, for example, in the medium class as**
18 **well ?**

19 The authors are thankful for pointing out the possible source of misunderstanding. With this
20 sentence we wanted to refer to the main source of underestimation which is represented by the
21 percentages indicated in Figure 6 (in the original manuscript). The overlaps mainly occur
22 between the low and medium susceptibility class (PP→ULCI 6%; LLCI→PP 5%) compared
23 to the percentage of overlapping cells of the medium and high susceptibility class (PP→ULCI
24 2%; LLCI→PP 2%).

25 Added to page 1029, line 27

26 Special attention should be given to the low susceptibility class. Here, the highest
27 percentage of overlapping classes and underestimation of the susceptibility were
28 detected.

1 30. *p1030 l3ff, especially l9ff: I feel that the comparison of your susceptibility map with*
2 *the hazard zonation plans is not fully feasible. Risks are induced by mass movements not*
3 *exclusively where they initiate, but also where they stop (in case of mass movements with a*
4 *considerable runout). For some types of movement, the hazard zonation map needs to*
5 *assess the runout zone as well. This might or might not be the case for earth and debris*
6 *slides that represent the main focus of this paper.]*

7 The authors agree with the expressed reservations of the referee. The general limitation of
8 susceptibility maps is the input data used for the modelling. In this study the landslide
9 inventory consisted of a point inventory representing the main scarps. Therefore, the modelled
10 susceptibility is also only showing the probability of the presence of main scarps. The runout
11 zone is not considered explicitly as no runout modelling was performed. However, it might
12 happen more or less accidentally that the possible runout zone might be located in the same
13 susceptibility class as the main scarp, due to the adjustment of the probability by the sampling
14 ratio and the later on classification in three classes only. We did an analysis of the coverage of
15 in the classified landslide susceptibility map by landslide polygons available in some parts of
16 the Flysch Zone. We found the majority of the landslide polygons (65%) were located in the
17 high susceptibility class. Furthermore, 29% were located in the medium and only 6% in the
18 low class. Naturally, this is dependent on the selected classification method or thresholds.

19 Accordingly, a sentence was included in the discussion section of the revised manuscript. The
20 discussion section 7.3 was changed significantly as we put the focus more on the need of
21 visualizing and communicating the results than on the comparison with the hazard zonation
22 plans. This comparison was removed from the revised manuscript. The changes of the
23 discussion in section 7.3 is presented here:

24 Added to page 1029, line 25

25 Some model form uncertainties within this method arise from using the lookup table
26 for transferring the prediction standard error to all grid cells as shown by the range of
27 resulting R^2 . This method might be improved or substituted by a function assigning the
28 standard errors to all grid cells.

29 It was found that in the classified map the majority of grid cells did not change.
30 However, there are differences between the modelling domains where some domains
31 had larger overlaps of different susceptibility classes than others. Special attention

1 should be given to the low susceptibility class. Here, the highest percentage of
2 overlapping classes and underestimation of the susceptibility were detected.

3 The visualization of these spatially varying uncertainties is of special interest for
4 future land-use and development planning usually performed by non-landslide experts.
5 In the aftermath of this study each landslide susceptibility class will be related to, not
6 legally binding, recommendations for the designation of new building areas.
7 Therefore, a misclassification (e.g. low instead of medium susceptibility) might lead to
8 an interpretation by the municipality or landowner that underestimated landslide
9 susceptibility. Knowledge about the susceptibility class overlaps might outline where
10 more caution and detailed investigations are necessary. Additionally, it also shows
11 where no uncertainties are expected, which might help to avoid costs for slope
12 investigations.

13 There is also a need to communicate the research results and their quality with
14 appropriate explanations for the local officials, environmental managers and the public
15 to raise awareness and knowledge on it which leads to an easier understanding and
16 incorporation of the results into the decision-making process (Knuepfer and Petersen,
17 2002; Rogers, 2006; Brierley, 2009; Hill et al., 2013). This analysis might aid to a
18 good acceptance of the landslide susceptibility maps in the local governments, as
19 instead of a fuzzy statement on involved uncertainties these are clearly shown in a map
20 on grid cell level (Guzzetti et al., 2006; Luoto et al., 2010). Furthermore, the
21 preparation of the susceptibility maps showing the class overlaps contributes to an
22 easier understanding of the possible effects of the prediction uncertainties.

23 The question if the policy makers or stakeholders are really interested in knowing
24 more about the uncertainty is discussed conversely. The study of Brugnach et al.
25 (2006) pointed out that the confidence in modelling results is dependent on the way
26 the uncertainties are addressed. Policy makers were missing more information on the
27 uncertainty of any model result. Therefore, the modelling results should be presented
28 with a measure of uncertainty or confidence indicator (Brugnach et al., 2006). In
29 habitat suitability modelling the visualisation of uncertainty was identified as relevant
30 to inform decision-makers about areas with extreme error, but also about areas which
31 are particularly well modelled (Elith et al., 2002). This openly addresses the
32 uncertainties involved in the maps instead of giving an impression of certainty (Elith

1 et al, 2002). However, interviews of Klimeš and Blahůt (2012) showed that local
2 governments do not want any information on uncertainties.

3 Nevertheless, these uncertainties might have severe consequences on buildings and
4 their inhabitants if an event occurred within the uncertainties of the method used to
5 delineate the hazard zones. The converse discussion shows, that more or better
6 communication with the stakeholders or policy makers (also during the modelling
7 process) is necessary to learn about uncertainties and enlarge confidence into the
8 modelling (Brugnach et al., 2006). However, the way how the uncertainties are
9 presented to the stakeholder has to be adapted by the scientist to ensure the success of
10 the communication. The visualization of some aspects of the quality of landslide
11 susceptibility maps, such as the spatially varying prediction uncertainty, can enhance
12 the communication among experts and decision-makers to facilitate informed
13 decisions (Kunz et al., 2011).

14 Additionally, further aspects of considering and communicating the effects of
15 epistemic uncertainty are still open research fields in susceptibility modelling. A clear
16 assessment of these is necessary to evaluate on their consequences on the
17 susceptibility (or hazard or risk) map.

18 **31. p1047 Fig. 4: Why does the AUROC for domain 230 scale on a 0-to-1 axis, while the**
19 **AUROC has a range of [0.5,1]?? This is one of two inconsistent uses of AUROC range (see**
20 **comment p1021 l13). ...**

21 The authors are thankful for pointing out the inconsistencies in the original manuscript. We
22 apologize for the inaccuracy and inserted some details on the AUROC. Please refer to our
23 reply on comment 20 for the changes made in the revised version of the manuscript.

24 **32. ... Moreover, the legend for each boxplot should be changed: either use "spCV and**
25 **nspCV" or (shorter and probably better) "sp and nsp".**

26 The legend of each boxplot was changed according to the suggestion of the referee.

1 33. *p1048 Fig. 5b: Something is wrong with the y axis labels. Either it should be the*
2 *numbers from 0.00 to 0.10, or from 0.00 to 1.00.*

3 We are grateful for the careful check of our figures. The y axis labels were adjusted to range
4 from 0.00 to 0.10.

5

6 **Technical Corrections**

7 34. *p1013 l13ff: I suggest to split this sentence: "Among the currently available*
8 *methods for landslide susceptibility modelling a GAM shows a compromise between the*
9 *flexibility of machine learning algorithms and the smooth representation which results*
10 *from GLMs such as logistic regression; meanwhile, it still gives the opportunity of a*
11 *transparent and easy interpretable model (Brenning, 2008; Goetz et al., 2011).*
12 *Alternatively: ...such as logistic regression while still giving the opportunity..."*

13 The sentence was changes according to the suggestion:

14 Added to page 1013, line 13

15 Among the currently available methods for landslide susceptibility modelling a
16 generalized additive model (GAM) is a compromise between the nonlinear predictive
17 flexibility of machine learning algorithms and the smooth, yet linear, predictions of
18 GLMs such as logistic regression. The model fit of the GAM can be easily
19 interpretable unlike most machine learning algorithms (Brenning, 2008; Goetz et al.,
20 2011).

21

22 **References**

23 Beguería, S.: Validation and Evaluation of Predictive Models in Hazard Assessment and Risk
24 Management, Nat Hazards, 37(3), 315–329, 2006.

25 Bell, R.: Lokale und regionale Gefahren- und Risikoanalyse gravitativer Massenbewegungen
26 an der Schwäbischen Alb, Rheinische Friedrich-Wilhelms-Universität Bonn, Bonn., 2007.

27 Beven, K. J. and Kirkby, M. J.: A physically based, variable contributing area model of basin
28 hydrology, Hydrological Sciences Bulletin, 24(1), 43–69, 1979.

- 1 Blume, H. P., Horn, R., Brümmer, G. W., Kandeler, E., Kögel-Knabner, I., Kretschmar, R.,
2 Stahr, K., Schachtschabel, P. and Wilke, B. M.: Scheffer/Schachtschabel: Lehrbuch der
3 Bodenkunde, Spektrum Akademischer Verlag.
- 4 Brenning, A.: Spatial prediction models for landslide hazards: review, comparison and
5 evaluation, *Natural Hazards and Earth System Sciences*, 5(6), 853–862, 2005.
- 6 Brenning, A.: Statistical Geocomputing combining R and SAGA: The Example of Landslide
7 susceptibility Analysis with generalized additive Models, *SAGA—Seconds Out*, 19, 23–32,
8 2008.
- 9 Brenning, A.: Improved spatial analysis and prediction of landslide susceptibility: Practical
10 recommendations, in *Landslides and Engineered Slopes, Protecting Society through Improved
11 Understanding*, edited by: Eberhardt, E., Froese, C., Turner, A. K., and Leroueil, S., Taylor &
12 Francis, Banff, Alberta, Canada., 789-795, 2012a.
- 13 Brenning, A.: Spatial cross-validation and bootstrap for the assessment of prediction rules in
14 remote sensing: the R package “sperrorest”, *IEEE International Geoscience and Remote
15 Sensing Symposium (IGARSS)*, 23-27 July 2012, Munich, 5372–5375, 2012b.
- 16 Breslow, N. E. and Day, N. E.: *Statistical methods in cancer research*, Iarc Lyon, France.
17 [online] Available from: [http://com.iarc.fr/en/publications/pdfs-online/stat/sp32/SP32_vol1-](http://com.iarc.fr/en/publications/pdfs-online/stat/sp32/SP32_vol1-0.pdf)
18 [0.pdf](http://com.iarc.fr/en/publications/pdfs-online/stat/sp32/SP32_vol1-0.pdf) (Accessed 16 July 2012), 1980.
- 19 Brierley, G.: Communicating Geomorphology, *J. of Geography in Higher Educ.*, 33(1), 3–17,
20 2009.
- 21 Brugnach, M., Tagg, A., Keil, F. and Lange, W. J.: Uncertainty Matters: Computer Models at
22 the Science–Policy Interface, *Water Resources Management*, 21(7), 1075–1090, 2006.
- 23 Chung, C. J. . and Fabbri, A. G.: Validation of spatial prediction models for landslide hazard
24 mapping, *Natural Hazards*, 30(3), 451–472, 2003.
- 25 Chung, C. J. F. and Fabbri, A. G.: Probabilistic prediction models for landslide hazard
26 mapping, *Photogrammetric Engineering and Remote Sensing*, 65(12), 1389–1399, 1999.
- 27 Demoulin, A. and Chung, C.-J. F.: Mapping landslide susceptibility from small datasets: A
28 case study in the Pays de Herve (E Belgium), *Geomorphology*, 89(3-4), 391–404, 2007.
- 29 Elith, J., Burgman, M. A. and Regan, H. M.: Mapping epistemic uncertainties and vague
30 concepts in predictions of species distribution, *Ecological modelling*, 157(2), 313–329, 2002.

- 1 Fabbri, A. G., Chung, C. J. F., Cendrero, A. and Remondo, J.: Is prediction of future
2 landslides possible with a GIS?, *Natural Hazards*, 30(3), 487–503, 2003.
- 3 Frattini, P., Crosta, G. and Carrara, A.: Techniques for evaluating the performance of
4 landslide susceptibility models, *Engineering geology*, 111(1-4), 62–72, 2010.
- 5 Goetz, J. N., Guthrie, R. H. and Brenning, A.: Integrating physical and empirical landslide
6 susceptibility models using generalized additive models, *Geomorphology*, 129, 376–386,
7 2011.
- 8 Guzzetti, F., Reichenbach, P., Ardizzone, F., Cardinali, M. and Galli, M.: Estimating the
9 quality of landslide susceptibility models, *Geomorphology*, 81(1-2), 166–184, 2006.
- 10 Hammerl, C. and Lenhardt, W.: *Erdbeben in Österreich*, Leykam, Graz, 1997.
- 11 Hanley, J.A. and McNeil, B.J.: The meaning and use of the area under a receiver operating
12 characteristic (ROC) curve, *Radiology*, 143, 29-36, 1982.
- 13 Heckmann, T., Gegg, K., Gegg, A. and Becht, M.: Sample size matters: investigating the
14 effect of sample size on a logistic regression debris flow susceptibility model, *Natural
15 Hazards and Earth System Sciences Discussions*, 1(3), 2731–2779, 2013.
- 16 Hill, L. J., Sparks, R. S. J. and Rougier, J. C.: Risk assessment and uncertainty in natural
17 hazards, in *Risk and uncertainty assessment for natural hazards*, edited by J. C. Rougier, R. S.
18 J. Sparks, and L. J. Hill, pp. 1–18, Cambridge University Press, Cambridge, 2013.
- 19 Hosmer, D. W. and Lemeshow, S.: *Applied logistic regression*, Wiley, New York, NY, 2000.
- 20 King, G. and Zeng, L.: Logistic regression in rare events data, *Political analysis*, 9(2), 137–
21 163, 2001.
- 22 Klimeš, J. and Blahůt, J.: Landslide risk analysis and its application in regional planning: an
23 example from the highlands of the Outer Western Carpathians, Czech Republic, *Natural
24 Hazards*, 2012.
- 25 Knuepfer, P. L. . and Petersen, J. F.: Geomorphology in the public eye: policy issues,
26 education, and the public, *Geomorphology*, 47(2–4), 95–105, 2002.
- 27 Kunz, M., Grêt-Regamey, A. and Hurni, L.: Visualization of uncertainty in natural hazards
28 assessments using an interactive cartographic information system, *Natural Hazards*, 59(3),
29 1735–1751, 2011.

- 1 Luoto, M., Marmion, M. and Hjort, J.: Assessing spatial uncertainty in predictive
2 geomorphological mapping: A multi-modelling approach, *Computers & Geosciences*, 36(3),
3 355–361, 2010.
- 4 Malamud, B. D., Turcotte, D. L., Guzzetti, F. and Reichenbach, P.: Landslide inventories and
5 their statistical properties, *Earth Surf. Process. Landforms*, 29(6), 687–711, 2004.
- 6 Remondo, J., González, A., De Terán, J. R. D., Cendrero, A., Fabbri, A. and Chung, C. J. F.:
7 Validation of landslide susceptibility maps; examples and applications from a case study in
8 Northern Spain, *Natural Hazards*, 30(3), 437–449, 2003.
- 9 Rogers, K. H.: The real river management challenge: integrating scientists, stakeholders and
10 service agencies, *River Research and Applications*, 22(2), 269–280, 2006.
- 11 Rossi, M., Guzzetti, F., Reichenbach, P., Mondini, A. C. and Peruccacci, S.: Optimal
12 landslide susceptibility zonation based on multiple forecasts, *Geomorphology*, 114(3), 129–
13 142, 2010.
- 14 Schwenk, H.: Massenbewegungen in Niederösterreich 1953 - 1990, in *Jahrbuch der*
15 *Geologischen Bundesanstalt*, vol. 135, pp. 597–660, Geologische Bundesanstalt, Wien., 1992.
- 16 Seibert, J., Stendahl, J. and Sørensen, R.: Topographical influences on soil properties in
17 boreal forests, *Geoderma*, 141(1-2), 139–148, 2007.
- 18 Trigila, A., Frattini, P., Casagli, N., Catani, F., Crosta, G., Esposito, C., Iadanza, C.,
19 Lagomarsino, D., Mugnozza, G., Segoni, S., Spizzichino, D., Tofani, V. and Lari, S.:
20 Landslide Susceptibility Mapping at National Scale: The Italian Case Study, in *Landslide*
21 *Science and Practice*, edited by C. Margottini, P. Canuti, and K. Sassa, pp. 287–295, Springer
22 Berlin Heidelberg. 2013.
- 23 Ury, H. K.: Efficiency of case-control studies with multiple controls per case: continuous or
24 dichotomous data, *Biometrics*, (31), 643–649, 1975.
- 25 Wacholder, S., Silverman, D. T., McLaughlin, J. K. and Mandel, J. S.: Selection of controls in
26 case-control studies: III. Design options, *American Journal of Epidemiology*, 135(9), 1042–
27 1050, 1992.
- 28 Van Westen, C. J., Castellanos, E. and Kuriakose, S. L.: Spatial data for landslide
29 susceptibility, hazard, and vulnerability assessment: An overview, *Engineering Geology*,
30 102(3-4), 112–131, 2008.