(This is a sample cover image for this issue. The actual cover is not yet available at this time.)

Contents lists available at SciVerse ScienceDirect

# Catena

# Plausibility test of conceptual soil maps using relief parameters

Markus Möller [a,*], Thomas Koschitzki [b], Klaus-Jörg Hartmann [c], Reinhold Jahn [d]

[a] University Halle-Wittenberg, Department of Remote Sensing and Cartography, Von-Seckendorff-Platz 4, 06120 Halle (Saale), Germany
[b] Geoflux GbR, Lessingstr. 37, 06114 Halle (Saale), Germany
[c] State Institute for Geology and Natural Resources Saxony-Anhalt, Köthener Str. 38, 06118 Halle (Saale), Germany
[d] University Halle-Wittenberg, Department of Soil Science and Soil Conservation, Von-Seckendorff-Platz 3, 06120 Halle (Saale), Germany

## ARTICLE INFO

## ABSTRACT

The motivation for this article results from the fact that conceptual soil maps show oftentimes inaccuracies with regard to soil unit boundaries or misfits between original paper and actual soil-related information. Using the example of a German conceptual soil map (CSM), we introduce a procedure which could be considered as a framework for testing the terrain-related plausibility applied within a genetic based soil-ordering system. Framework means that all tests and the underlying methods can be adapted to specific targets. The procedure enables both reproducible integration of expert knowledge and application of statistically sound methods.

The CSM of the German Federal State of Saxony-Anhalt was tested regarding the plausibility of colluvial and fluvial process domains. The plausibility test consists of four steps and was exemplified on a study area of 100 km². First, basic relief parameters were combined to the explaining relief parameters *Floodplain Index (FPI)* and *Mass Balance Index (MBI)* enabling a classification of process domains by relative descriptions. Second, relief parameters and aggregated CSM soil units were integrated to soil-terrain objects (STO) executing a region-growing segmentation algorithm. In the third step, the one-dimensional *MBI* or *FPI* feature space of STO entities were clustered by using the K-means algorithm. The fourth step comprises the expert-based selection of reference clusters (RC) representing colluvial and fluvial process domains accepted as being true. Then, empirical cumulative distribution functions (ECDF) of RC and remaining soil unit-related STO clusters were compared by a traditional goodness-of-fit test whose suitability for estimation of terrain-related CSM plausibility is shown. Finally, the resulting ECDF distances were visualized.

The testing procedure could also be used for the supervised selection of appropriate samples for automatic classification algorithms. The data integration approach is generally suitable for adopting existing data in computer-based systems.

## 1. Introduction

Conceptual soil maps (CSM) are the result of an expert-based integration process wherein different soil-related information or already existing older soil maps are combined by soil surveyors (Dobos and Hengl, 2009). The resulting soil maps "are representations of structured knowledge" (Bui, 2004). The data integration process is mainly guided by the ordering system. Genetic soil-ordering concepts are based on generally accepted perceptions of soil genesis and allow more expert-based interpretation than classification systems where threshold defined diagnostic horizons, features and the horizon sequences determine soil units (Albrecht et al., 2005; Buol et al., 2003).

A German CSM example is the preliminary soil map 1:50,000 of Saxony-Anhalt (in German: "Vorläufige Bodenkarte 1:50,000" or VBK 50; Hartmann, 2005, 2006). The VBK 50 map results from an expert-based integration process where older soil maps were semantically transferred into the actual (genetic) German soil ordering system (Ad-hoc AG, 2005). The soil map contains typical CSM inaccuracies. First, misfits between original paper and actual, more accurate soil-related information exist. Second, locations of systematic soil unit boundaries are often incorrect due to their subjective delineation. This can be observed especially between areas of fluvial and terrestrial process domains. A special problem is related to the attribute structure of the used older soil maps which were the basis of the CSM creation (Deumlich et al., 1998; Müller and Volk, 2001). The attributes describe heterogeneous soil units of genetically linked soils. That means that the occurrence of some – especially colluvial – soil units are only listed in the attribute table but not represent polygons (Möller, 2008).

The mentioned inaccuracies are mostly terrain-related and concern especially colluvial and fluvial process domains. Information about surface topography can nowadays be derived from easily accessible digital elevation models (DEM) in different spatial resolutions and accuracies (Hengl and MacMillan, 2009). Thus, the main objective of this study is

---

* Corresponding author. Tel.: +49 345 5526025; fax: +49 345 2394019.
*E-mail address:* markus.moeller@geo.uni-halle.de (M. Möller).
*URL:* http://www.geo.uni-halle.de/geofern/mitglieder/moeller/ (M. Möller).

the development and application of a procedure testing the colluvial and fluvial plausibility of VBK 50. Furthermore, heterogeneous soil units should be geometrically disaggregated regarding the occurrence of colluvial process domains.

In this study, the plausibility test is demonstrated on the example of the German topographic map TK25 4336 Könnern at a scale of 1:25,000 with an area of about 100 km². In a joint project of a state authority, a scientific institution and an engineering office, the procedure was applied on the total area of the German Federal State of Saxony-Anhalt ($\approx$20,000 km²). The project's outcome can be considered as a compromise solution, in which pedological and digital soil mapping (DSM) expertise were balanced.

We are using the term *plausibility* instead of *quality* or *accuracy*. Quality and accuracy are related to international standards for geodata (e.g. ISO 19138, 2006). This should help data producers objectively describe the quality of data and determine its quality using statistical calculation rather than subjective estimation. Although we support the utilization of statistical quality measures, we also recognize the need for possibilities to deal with expert knowledge in a reproducible manner (see Deumlich et al., 2010). The integrating result of subjective and objective quality to assess the soil map's goodness we refer to as *plausibility*.

In this study, plausibility is considered as distance between reference and test distributions of explaining relief parameters. We show how basic relief parameters can be combined to specific indices which explain the occurrence of colluvial and fluvial process domains. Furthermore, we demonstrate how reference distributions can be defined by expert-knowledge in a transparent and traceable manner. Finally, the suitability of a traditional goodness-of-fit-test for comparing relief parameters' distributions is investigated.

## 2. Materials and methods

### 2.1. Study site

The study area is situated in the German Federal State of Saxony-Anhalt and represents heterogeneous soil and relief conditions (Fig. 1, Table 1; Möller et al., 2008). The formation of parent materials, relief and soil formation was connected with glacial and periglacial conditions during the Saalian and Weichselian glacial stages where plateaus and floodplains were shaped. Plateaus and plateau margins are mainly covered by Weichselian loess and Saalian moraine material. There, calcareous Ah/C and black soils dominate (Pararendzina, Tschernosem). Where older sandstones, clay or limestones of mainly
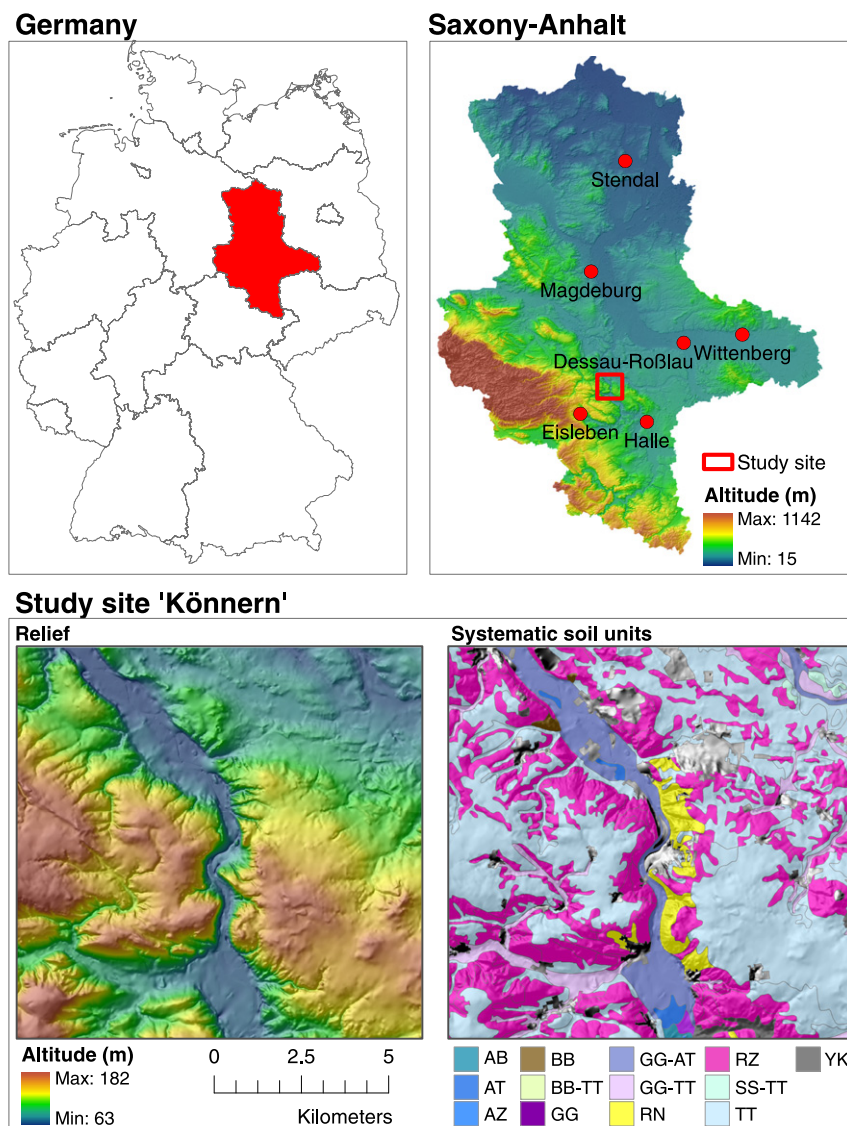


**Fig. 1.** Study site location as well as its soil (see Table 1) and terrain conditions (Data sources: Soil units — http://www.lagb.sachsen-anhalt.de | DEM — http://www.lvermgeo.sachsen-anhalt.de).

**Table 1**
German systematical soil units according to Ad-hoc AG (2005) and the most probable reference soil group according to IUSS Working Group (2006).

| | German soil unit | Short description | Most probable reference soil group |
|---|---|---|---|
| RN | Ranker | Ah/C soils from silicatic rock | Haplic Leptosol from periglacial layers of sandstone |
| RZ | Pararendzina | Calcareous Ah/C soils | Haplic Regosol Calcaric from periglacial layers of loess or marly materials |
| TT | Tschernosem | Black soils | Haplic Chernozem from loess |
| BB | Braunerde | Brown soils | Haplic Cambisol from periglacial layers of sandstone |
| YK | Kolluvisol | Soils from eroded top soil material | Colluvic Regosol from colluvic loess material |
| AB | Vega | Brown floodplain soils | Endofluvic Cambisol from fluvic material |
| AT | Tschernitza | Chernozem like floodplain soils | Endofluvic Chernozem from fluvic material |
| AT | Kalkpaternia | Calcareous Ah/C floodplain soils | Haplic Fluvisol from fluvic material |
| GG | Gley | Groundwater affected soils | Endofluvic Endogleyic Cambisol from fluvic material |

Permian and Carboniferous ages emerge, brown soils and Ah/C soil from silicatic rock occur (Braunerde and Ranker). Floodplain and groundwater affected soils have developed in the fluvial sediments (Vega, Tschernitza, Kalkpaternia, Gley). The occurrence of colluvial soils (Kolluvisol) is favored by the intensive agriculture and intense summer rainstorm events.

## 2.2. Workflow

The procedure consists of four steps (Fig. 2). On the basis of a digital elevation model (DEM), the starting point of the workflow is the
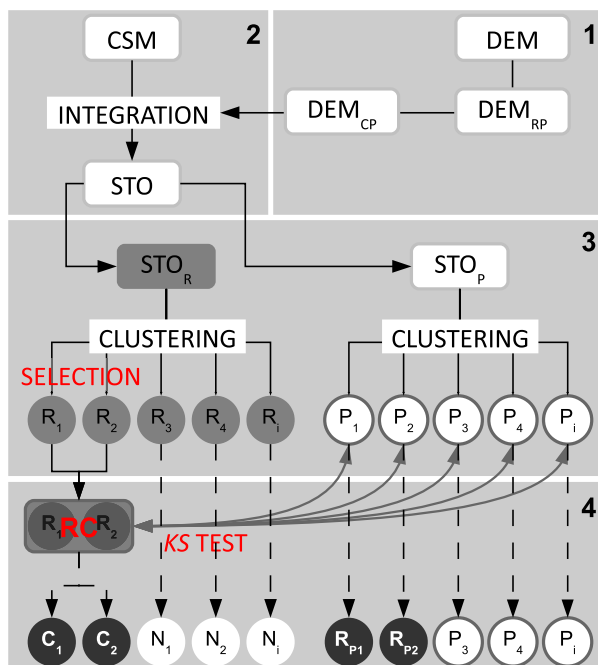


**Fig. 2.** Workflow. DEM — digital elevation model | $DEM_{RP}$ — basic relief parameter | $DEM_{CP}$ — combined relief parameter on demand | CSM — conceptual soil map | STO — soil-terrain object | $STO_P$ — STOs which should be tested | $STO_R$ — STOs representing the process domain of interest based on CSM | RC — $STO_R$ clusters representing terrain-related process domains of interest selected by experts | $R_i$ — $STO_R$ clusters | $P_i$ — $STO_P$ clusters | $C_i$ — $R_i$ which are confirmed regarding terrain-related process domains | $N_i$ — $R_i$ which are rejected regarding terrain-related process domains and have to be newly classified | $R_{Pi}$ – $STO_P$ clusters which are tested positive regarding terrain-related process domains.

derivation of appropriate relief parameters (Section 2.2.1). Our basic approach is to combine several basic relief parameters ($DEM_{RP}$) to as few as possible combined relief parameters on demand ($DEM_{CP}$). The parameter combination is guided by the objective of the test. In this study, colluvial and fluvial process domains should be characterized. The combination is carried out in such a way that process domains can be classified by relative descriptions. Relative definitions like minimum or maximum values enable an easier integration of expert knowledge. In addition, classification models can be better transferred to other geomorphologic regions (Dragut and Blaschke, 2006; Möller et al., 2008).

The data integration procedure in the second step of the workflow aims at the coupling of the CSM data set and the predicting relief parameters ($DEM_{CP}$). For this purpose, the CSM spatial domains were subdivided into soil-terrain objects (STO) by the application of a segmentation algorithm (Section 2.2.2). STOs can be characterized as groups of pixels of different relief parameters which are aggregated to landform elements according to a scale-specific homogeneity (Dragut and Eisank, 2011; MacMillan and Shary, 2009; Minár and Evans, 2008) considering already existing soil map boundaries. In step 3, two STO groups have to be defined:

1. $STO_P$ are segmented systematical soil units which should be tested regarding their terrain-related plausibility. In this study, soil units were aggregated according to their dominating soil-terrain-related formation (Table 2). The areas of aggregated soil units are shown in Fig. 5a. Accordingly, the aggregated soil units T and R/B dominate spatially while the units A/G and especially Y are represented by smaller areas.
2. $STO_R$ stands for a segmented soil unit which represents the process domain of interest provided by the original CSM. In Table 2, the associated (here: colluvial and fluvial) CSM soil units are gray emphasized. While the soil unit YK is affected by colluvial processes, the soil units GG, AB, AT, AZ, GG-AT are mainly influenced by fluvial processes. In the following, both groups are classed as Y and A/G. $STO_{A/G}$ and $STO_Y$ are the basis for the selection of reference clusters (RC).

The clustering procedure can be considered as statistical grouping of a feature space (Section 2.2.3). The grouping refers to the specific $DEM_{CP}$ feature space of STO entities. Each STO entity covers an aggregated soil unit and is separately clustered. According to Table 2, four entities are to be clustered in this study ($STO_Y$, $STO_{A/G}$, $STO_T$, $STO_{R/B}$). The grouping leads to $STO_R$ and $STO_P$ clusters which are labeled as $R_i$ and $P_i$. Reference clusters (RC) are determined by experts. This crucial operation is supported by the specific $DEM_{CA}$ properties which enable a relative classification (see step 1). The consideration of $R_i$ cluster values, their visualization in feature space and as segmented and clustered soil map helps to identify RCs. In other words, RC are $STO_R$ (here: $STO_Y$ or $STO_{A/G}$) with a representative value distribution of the (here: colluvial or fluvial) process domain of interest which is accepted as being true.

**Table 2**
Aggregation criteria for soil units. The gray emphasized soil units represent fluvial (A/G) and colluvial process domains of interest (Y).

| Aggregated soil unit | Soil unit | Aggregation criterion |
|---|---|---|
| A/G | GG,* AB, AT, AZ, GG-AT | Dominating fluvial processes |
| R/B | BB, RN, RZ | Dominating processes of solifluction, humus layer thickness < 4 dm |
| T | BB-TT, GG-TT, SS-TT, TT | Dominating processes of solifluction, humus layer thickness > 4 dm |
| Y | YK | Dominating processes of solifluction and colluvial process domains |

* In fact, GG stands for soils with gleyic properties (Gleysols). However, all Gleysols of the study area are situated within floodplains. Thus, these soils are also characterized by fluvial properties.
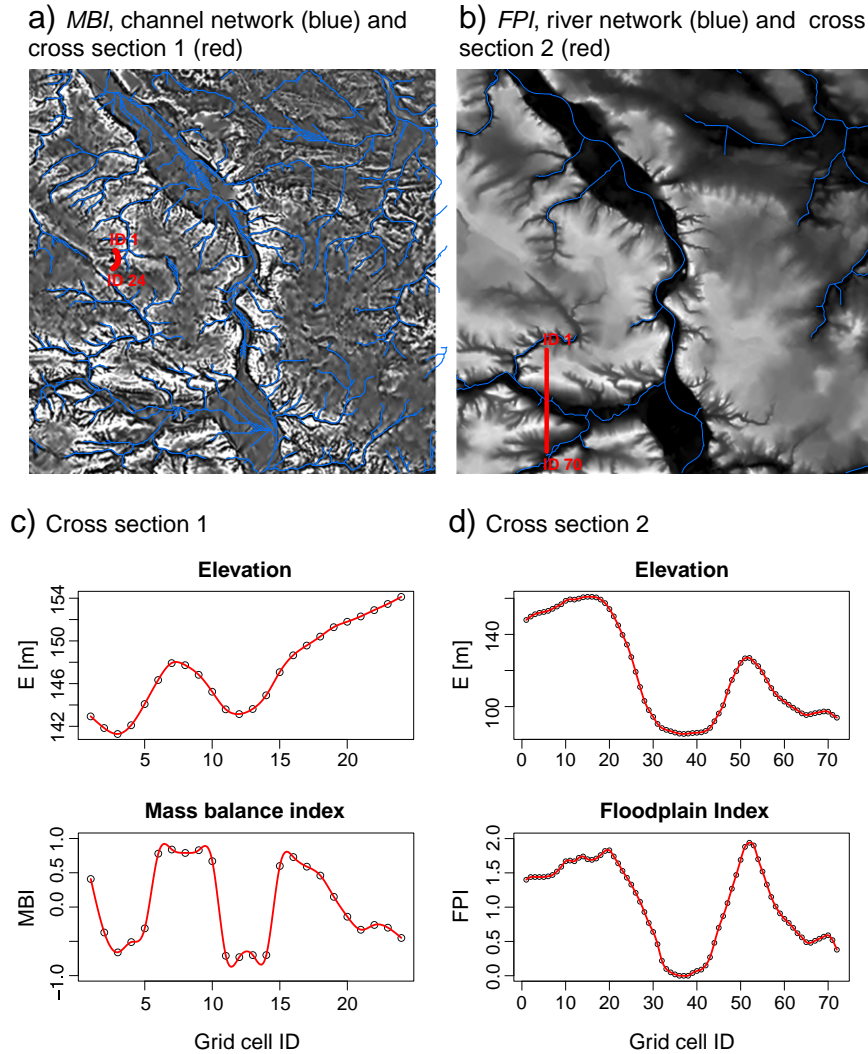
a) *MBI*, channel network (blue) and cross section 1 (red)

b) *FPI*, river network (blue) and cross section 2 (red)

c) Cross section 1

d) Cross section 2



**Fig. 3.** Relations between DEM cross sections and value ranges of *MBI* (a, c) and *FPI* (b, d).

Step 4 of the procedure comprises the comparison of RC and $P_i$ distributions by the application of a goodness-of-fit test (Section 2.2.4). As a result, four types of clusters emerge:

1. $STO_R$ clusters which are confirmed ($C_i$),
2. $STO_R$ clusters which are rejected and have to be newly classified ($N_i$),
3. $STO_P$ clusters which are statistically similar to RC ($R_{Pi}$) and
4. $STO_P$ clusters which are not similar to RC ($P_i$).

### 2.2.1. Digital elevation model and calculation of relief parameters

The used state-wide available digital elevation model (DEM) was originally generated by the digitalization of elevation contours of topographic maps in a scale of 1:10,000. The ANUDEM algorithm by Hutchinson (1989) was applied in order to create a hydrological sound DEM with a resolution of 20 m (see Fig. 1).

The calculation of basic relief parameters was performed within SAGA GIS and RSAGA environment using the application rsaga. geoprocessor[1] (Brenning, 2008; Olaya and Conrad, 2009). *Slope* (*S*) and *Total Curvature* (*TC*) are standard relief parameters and were calculated according to Zevenbergen and Thorne (1987). The parameter *Vertical Distance to Channel Network* (*VDN*) is the difference between the original elevation *DEM* and the interpolated channel network base

level (*DEM$_{BASE,N}$*; Eq. (1)). The calculation of the parameters *Vertical Distanceto River Network* (*VDR*) and *Vertical Distance to Culmination Network* (*VDC*) is similar to Eq. (1) however the reference for the base level interpolation differs. *VDR* uses the river network as reference (*DEM$_{BASE,R}$*, Eq. (2)). The comparison of Fig. 3a and b reveals the difference between river and channel network. *VDC* is referring to the culmination network resulting from a reversed DEM (*DEM$_r$*). Here, the base level is named as *DEM$_{r,BASE,C}$* (Eq. (3)).

$$VDN = DEM - DEM_{BASE,N} \tag{1}$$

$$VDR = DEM - DEM_{BASE,R} \tag{2}$$

$$VDC = DEM_r - DEM_{r,BASE,C} \tag{3}$$

In addition, the well-known *Topographic Wetness Index* (*TWI*) was calculated according to Eq. (4) where *A* is the *Specific Catchment* and *S* the *Slope* (Quinn et al., 1995).

$$TWI = ln\left[\frac{A}{tan(S)}\right] \tag{4}$$

All relief parameters were transformed into a unique value range by the application of Eq. (5). *x* stands for the corresponding relief parameter and *F* is a user-defined transfer constant affecting the

---

[1] http://cran.r-project.org/web/packages/RSAGA/RSAGA.pdf.

parameter's value distribution (Friedrich, 1996; Friedrich, 1998). Here, *F* values of 15 (*S, VDN, VDR, VDC, TWI*) and 0.01 (*TC*) were used leading to balanced ratio of dominating (e.g. floodplains) and smaller landforms (e.g. depressions; Möller et al., 2008).

$$f(x) = \frac{x}{|x| + F} \tag{5}$$

with $x = S, VDN, VDC, VDR, TWI, TC$; $f(S, VDN, VDC, VDR, TWI) \in [0, 1]$; $f(TC) \in [-1, 1]$

### 2.2.2. Segmentation

In this study, the region growing segmentation algorithm *Fractal Net Evolution Approach* (FNEA) described in detail by Baatz and Schäpe (2000) was applied to relief parameters. The algorithm relies on seed pixel groups with both the smallest (here: Euclidean) distance in pixel raster and in n-dimensional feature space of the used relief parameters. Then, the seeds grow as far as a halting criterion is fulfilled. The halting criterion could be a specific object heterogeneity or existing boundaries. The segmentation process leads to different aggregation levels of discrete landform elements. Each level represents a specific target scale consisting of objects with a comparable heterogeneity. The segmentation results can be influenced by parameters which allow the adaptation of the target segment's heterogeneity and shape (Möller et al., 2008).

FNEA has been proven as a suitable algorithm for detecting objects having meaning for soil-terrain-related issues (Dragut and Eisank, 2011). A crucial point is the determination of an optimal segmentation parameter setting (Dragut et al., 2009). Similar to Dragut and Blaschke (2006), we compared different segmentation results with significant known landforms like valleys and slope positions representing minimal object sizes.

### 2.2.3. Cluster analysis

Clustering belongs to the standard techniques of unsupervised learning and aims at the grouping of similar objects. In contrast to the aforementioned FNE algorithm, similarity only refers to feature space of data points (here: explaining relief parameters). The applied K-means algorithm uses the squared Euclidean distance as dissimilarity measure. The algorithm is described in detail by Hastie et al. (2009). Starting with a user-defined number of initial K centroids, each data point is iteratively assigned to the nearest cluster centroid. The maximum number of iterations must be specified by the analyst.

### 2.2.4. Comparing distributions

The comparison of reference clusters (RC) and clustered soil units ($P_i$) is done by the Kolmogorov Smirnov goodness-of-fit (KS) test (Davis, 2002; Thas, 2010). The main advantage of this nonparametric two-sample test – especially in dealing with environmental data – is the fact that any kind of distribution can be compared without requiring specific statistical conditions. Based on the empirical cumulative distribution function (ECDF) the KS test verifies whether two distributions are the same (null hypothesis) or significantly different from each other. The degree of difference is expressed by the maximal absolute difference *D* between the cumulative distributions of RC and $P_i$ (Eq. (6)). Both K-means clustering and KS test were executed within the statistical environment R[2] using the functions `kmeans`[3] and `ks.test`.[4]

$$D = max|ECDF_{RC} - ECDF_{Pi}| \tag{6}$$

---

[2] http://www.r-project.org.
[3] http://stat.ethz.ch/R-manual/R-devel/library/stats/html/kmeans.html.
[4] http://stat.ethz.ch/R-manual/R-devel/library/stats/html/ks.test.html.

## 3. Results

### 3.1. Step 1: Combination of relief parameters on demand

The *Mass Balance Index* (*MBI*) was used to characterize colluvial soil process domains. Negative *MBI* values represent areas of net deposition such as depressions and valleys, positive *MBI* values indicate areas of net erosion such as convex hill slopes, *MBI* values close to 0 refer to areas with a balance between erosion and deposition such as plain areas (Fig. 3a, c). *MBI* results from the combination of transformed relief parameters *f(S)*, *f(VDN)* and *f(TC)* (Eq. (7), see Section 2.2.1; Möller et al., 2008).

$$MBI = \begin{cases} f(TC) \times (1 - f(S)) \times (1 - f(VDN)) \text{ for } f(TC) < 0 \\ f(TC) \times (1 + f(S)) \times (1 + f(VDN)) \text{ for } f(TC) > 0 \end{cases} \tag{7}$$

As the name implies, the *Floodplain Index* (*FPI*) enables the detection of floodplains. Floodplains are located lower than their surroundings which can be expressed by the relief parameter *Relative Slope Position* (*RSP*). *RSP* is calculated according to Eq. (8) from the quotient of *f(VDR)* and the sum of *f(VDR)* and *f(VDC)* (see Section 2.2.1). Furthermore, floodplains are characterized by their maximal flatness or minimal slope and maximal flow accumulation which is represented by the relief parameter *TWI*. The combination of *RSP* and *f(TWI)* results in *FPI* (Eq. (9)). Floodplains can be detected by minimal *FPI* values regardless of their absolute altitude (Fig. 3b, d).

$$RSP = \frac{f(VDR)}{f(VDR) + f(VDC)} \tag{8}$$

$$FPI = \frac{RSP}{f(TWI)} \tag{9}$$

### 3.2. Step 2: Data integration

The segmentation operation was applied to the transformed relief parameters *f(TC)*, *f(S)*, *f(VDN)* and *f(VDR)*. The VBK 50 boundaries acted as additional halting criteria. The segmentation led to 34,749 soil-terrain objects (STO). This means that the pixel number of 324,220 was decreased to about one tenth.

The blue-framed objects in Fig. 4a represent STOs which subdivide a red-framed superior VBK 50 unit. The STOs vary in size depending on their terrain position. This corresponds to the soil unit-related distributions of STO sizes (Fig. 4b): While the units T and A/G occur in rather flat areas (green and orange-colored STOs), the units R/B and Y are connected to steeper areas (magenta and blue-colored STOs).

The data integration process results in soil unit-related *FPI* and *MBI* distributions (Fig. 5b and c). All medians show typical positions within *FPI* and *MBI* value ranges. This is particularly true for medians of $A/G_{FPI}$ and $Y_{MBI}$ which are minimal and confirm the assumptions in Section 3.1. However, all distributions are characterized by overlaps and wide value ranges.

### 3.3. Step 3: Cluster analysis and identification of reference clusters

As already stated in Section 2.2.3, K-means clustering requires determination of the desired cluster number. Here, the *FPI* and *MBI* feature space of each aggregated soil unit was subdivided into ten clusters, and twenty iterations were run for each clustering. In Fig. 6, the associated *FPI* and *MBI* box plots as well as their positions within *FPI* and *MBI* value ranges are visualized. The vertical red and blue solid lines correspond to the upper boundary of *MBI* and *FPI* value ranges which represent reference clusters (RC). RC can be considered as core zones of fluvial and colluvial process domains. They are composed of clusters which were identified by expert knowledge.
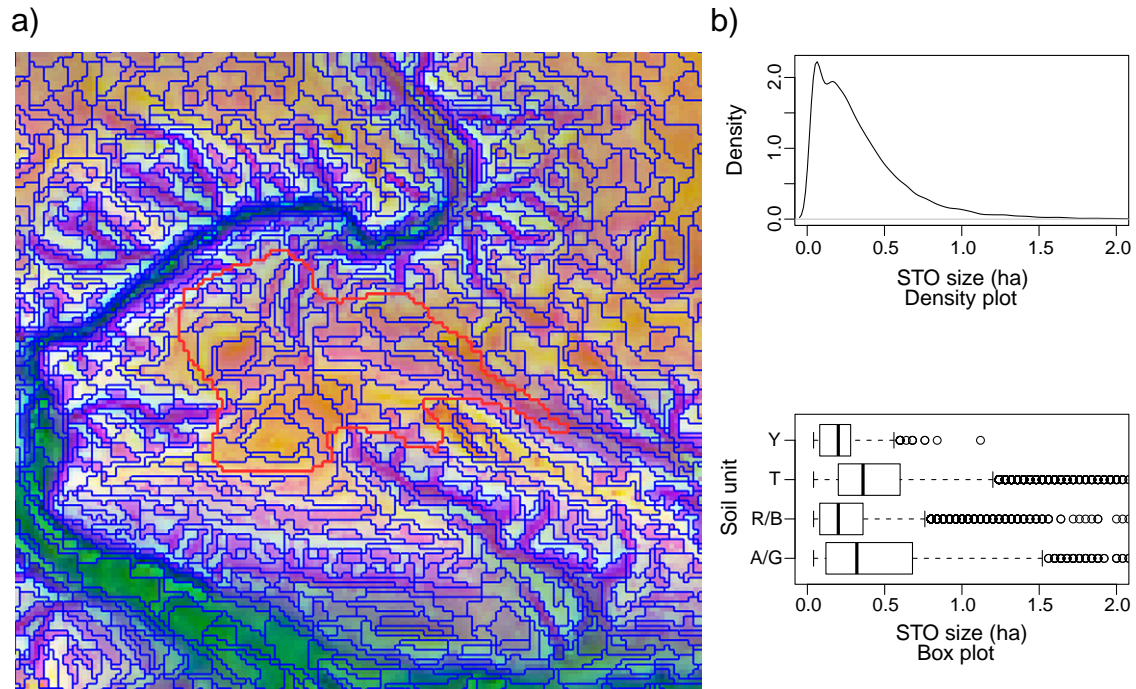
**Fig. 4.** Visualization of soil-terrain objects (STOs) on the example of a study site subset (a) and distributions of aggregated soil unit-related STO sizes (b).

RC for colluvial processes consists of the $MBI_Y$-clusters 3 and 6 (Fig. 6a). Fluvial process domains are described by the $FPI_{A/G}$-clusters 1, 4, 5, 6 and 9 (Fig. 6b). The gray dashed lines show the RC positions in the diagrams of the other soil units. In Fig. 8, the identified core zones correspond to areas with $D = 0$.

### 3.4. Step 4: Kolmogorov Smirnov (KS) goodness-of-fit-test

The KS test was carried out for both fluvial and colluvial processes. While the first test was applied to the aggregated soil units R/B, T and Y, the latter one was executed to the units R/B, T and A/G. In Fig. 7, for
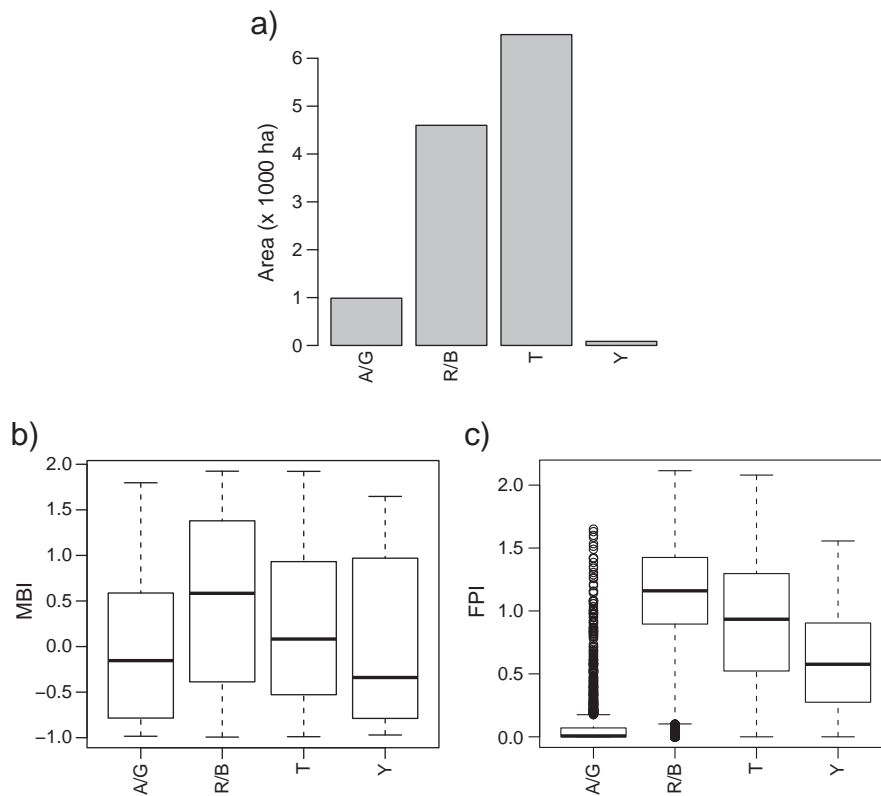


**Fig. 5.** Areas of aggregated soil units (a) as well as corresponding *MBI* (b) and *FPI* box plots (c).

## a) Identification of colluvial reference cluster (RC) within the aggregated soil unit Y



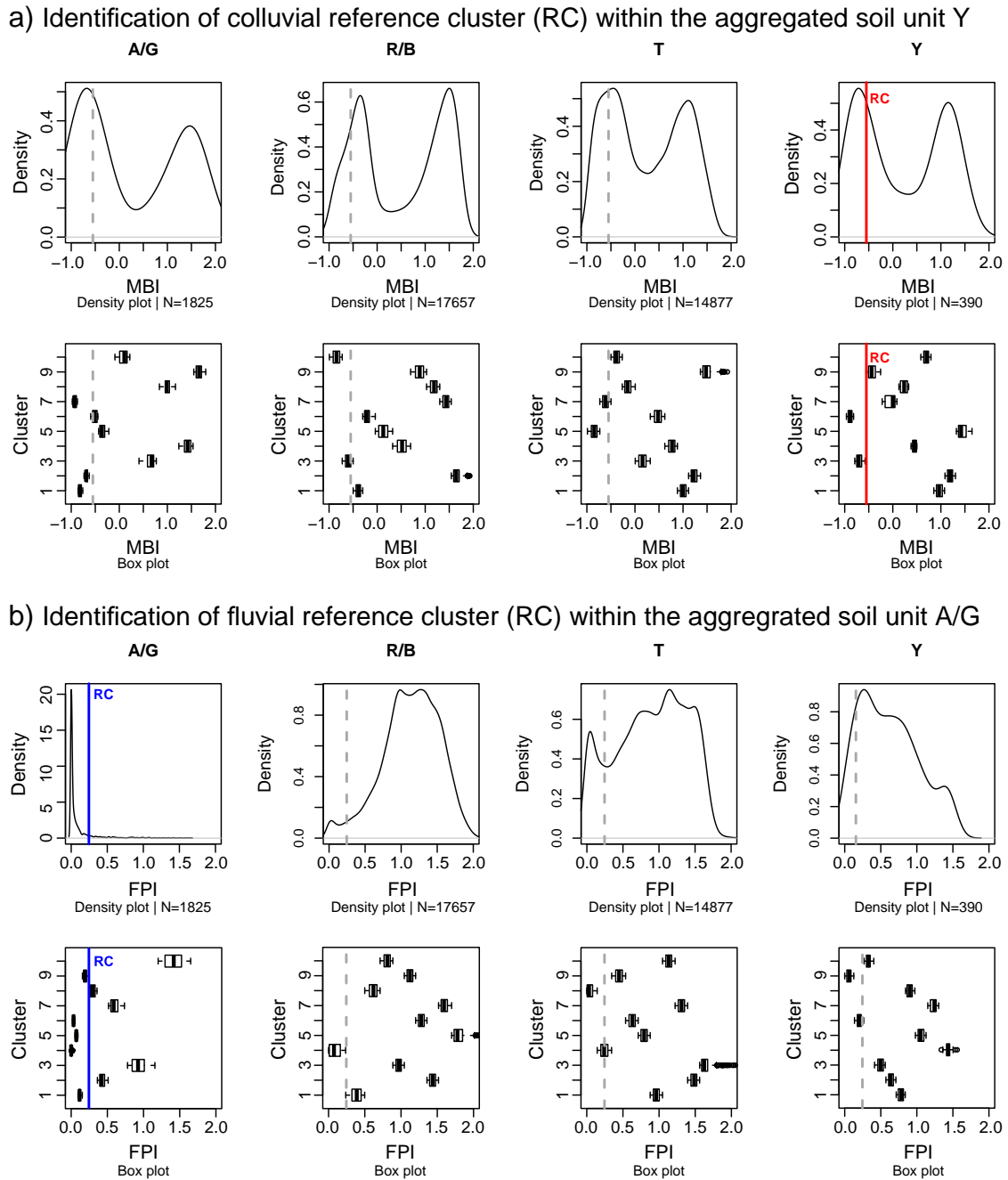## b) Identification of fluvial reference cluster (RC) within the aggregated soil unit A/G



**Fig. 6.** *MBI* and *FPI* density plots and cluster-related box plots for each aggregated soil unit. The vertical solid lines indicate reference clusters (RC) of colluvial (red) and fluvial process domains (blue). The gray dashed lines show the RC positions in the diagrams of the other soil units. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)
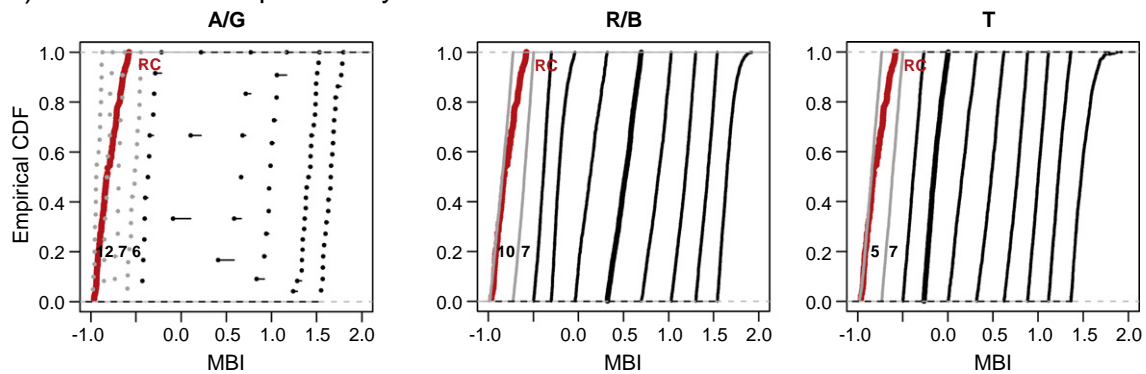
each aggregated soil unit, all cluster-related ECDFs are plotted as black graphs against the blue and red colored ECDF graphs of RC. ECDFs with differences $D<1$ are plotted as gray graphs. The associated cluster number is labeled and can be connected with the listed $D$ values of Table 3. Although all $p$ values denote significant differences between RC and cluster-related distributions, small $D$ values indicate similarities to colluvial or fluvial process domains.

In Fig. 8, the resulting $D$ values are joined with STOs and mapped. Four groups of STOs arise (see also Fig. 2 and Section 2.2):

1. Level $D=0$ denotes $STO_R$ which were identified as reference clusters (RC). These core zones of colluvial and fluvial process domains (class C) are displayed in dark blue.

2. $STO_R$ with $D=1$ were rejected as reference clusters. They are red highlighted and have to be newly classified (class $N_i$).

3. $STO_P$ with $D>0$ and $D<1$ indicate potential areas affected by fluvial or colluvial processes. With it, users can formulate a basis for decision-making which $STO_P$ should be newly assigned to fluvial or colluvialprocess domains (class $R_{Pi}$). In doing so, the application of a classification rule (here $D_Y \leq 0.38$, $D_{A/G} \leq 0.49$) would lead to a revised soil map with new soil unit-related *FPI* and *MBI* distributions (see Table 3).

The boxplot comparison of the original (Fig. 5b and c) and tested soil units demonstrates that value ranges of A/G and Y – which were identified as implausible – were removed and transferred into

## a) Test for colluvial plausibility
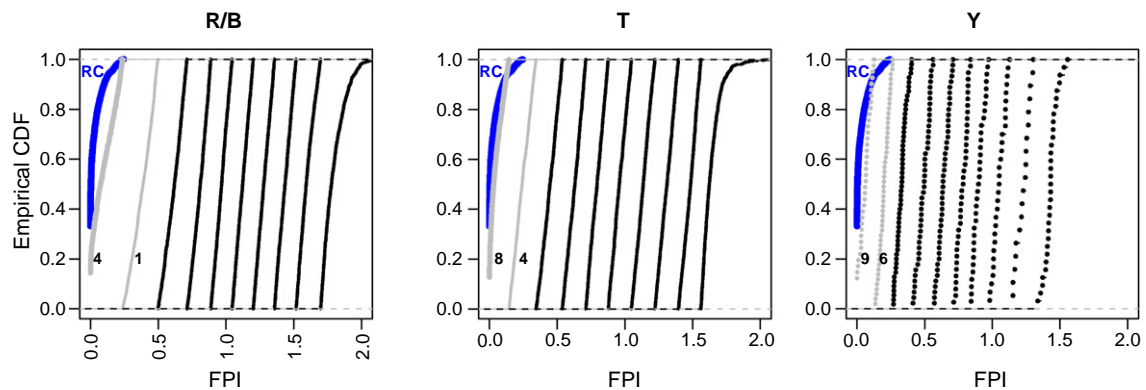


## b) Test for fluvial plausibility



**Fig. 7.** ECDF plot comparison of reference clusters (RC) and soil unit-related *MBI* and *FPI* clusters. Blue colored graphs indicate fluvial process domains, and red colored RC graphs stand for colluvial process domains. ECDF plots with *D*<1 are gray colored and labeled with the associated cluster number of Table 3. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

the new classes 'AG rejected' and 'Y rejected' (Fig. 9b and c). Both correspond to class N$_i$ of the workflow (Fig. 2). The classes A/G and Y consist of confirmed (C$_i$) and positive tested STO clusters (R$_{Pi}$). T and R/B correspond to negative tested clusters (P$_i$). Overlapping distributions between the soil units A/G or Y and R/B or T could be minimized. However, the remaining aggregated soil units A/G and Y still show overlapping distributions which were tested positive regarding both existence of fluvial and colluvial process domains. This is an expression of situations where both soil forming processes could take place and superimpose. This class is labeled as 'A/G or Y' and also corresponds to class R$_{Pi}$ of the workflow.

The comparison of original and revised A/G soil unit areas reveals an increase by almost half (Figs. 9a and 5a). The tenfold increase of

colluvial areas is caused by the specific thematic attribute structure of the older soil maps which formed the basis for the CSM creation (see Section 1). There, smaller colluvial soil units are only documented in the attribute table. This information has got lost during the semantic transformation. This means that the classification result represents a geometric disaggregation revealing semantic terrain-related information. The main source for the new assigned colluvial soils is the aggregated soil unit T (Tschernosem).

## 4. Discussion and conclusion

A main challenge of digital soil mapping (DSM) is the adoption of appropriate techniques and input data for operational use. A key to gain acceptance and overcome possible resistance is the development and application of standardized protocols for producing predictive maps (Hengl and MacMillan, 2009; MacMillan, 2008). Against this background, the presented procedure can be considered as a modular structured framework for terrain-related plausibility tests of conceptual soil maps (CSM; see Fig. 2). Framework means that all tests and the underlying methods can be adapted to specific test objectives. In this article, the CSM of the German Federal State of Saxony-Anhalt should be tested against their colluvial and fluvial plausibility. The applied procedure was exemplified on a test site with an area of about 100 km$^2$ and consists of four steps:

1. In the first step, explaining relief parameters on demand were calculated for the identification of fluvial and colluvial process domains. In doing so, basic relief parameters were combined to the *Floodplain Index* (*FPI*) and *Mass Balance Index* (*MBI*) allowing a relative characterizing of process domains. As shown by Möller et al. (2008), the

**Table 3**
D values between reference clusters (RC) and soil unit-related *MBI* and *FPI* clusters (see Fig. 7). Values of *D*<1 are gray and values of *D*<0.5 are bold emphasized.

| Aggregated soilunit | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | MBI clusters | | | | | |
| A/G | 0.38* | 0.68* | 1* | 1* | 1* | 0.96* | 0.67* | 1* | 1* | 1* |
| R/B | 1* | 1* | 0.72* | 1* | 1* | 1* | 1* | | 1* | 0.31* |
| T | 1* | 1* | 1* | 1* | 0.32* | 1* | 0.69* | 1* | 1* | 1* |
| | | | | | FPIclusters | | | | | |
| R/B | 0.99* | 1* | 1* | 0.39* | 1* | 1* | 1* | 1* | 1* | 1* |
| T | 1* | 1* | 1* | 0.95* | 1* | 1* | 1* | 0.30* | 1* | 1* |
| Y | 1* | 1* | 1* | 1* | 1* | 0.94* | 1* | 1* | 0.49* | 1* |

* Significance level *p*<0.01.

a) Test for colluvial plausibility

b) Test for fluvial plausibility



| D levels | 0 | >0 - 0.32 | >0.32 - 0.38 |
| | >0.38 - 0.72 | >0.72 - 0.99 | 1 |

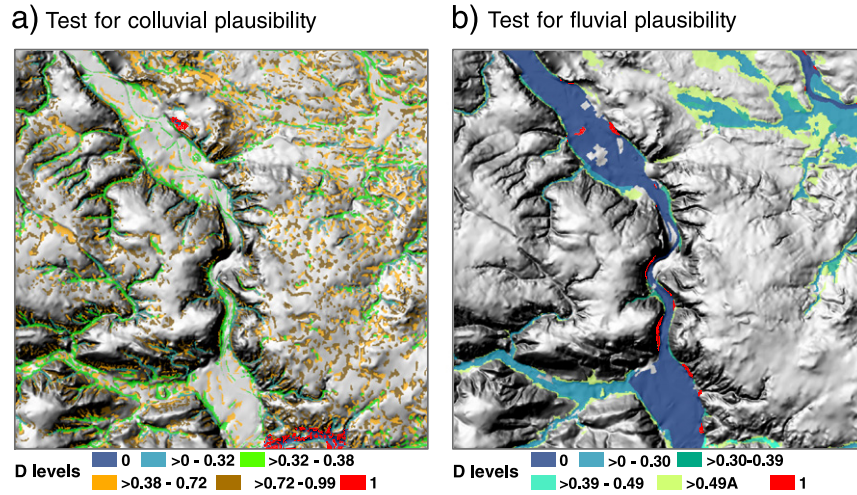| D levels | 0 | >0 - 0.30 | >0.30-0.39 |
| | >0.39 - 0.49 | >0.49A | 1 |

**Fig. 8.** Visualization of *D* value levels resulting from plausibility tests (see Table 3).

*MBI* value distribution can be affected by changing the transfer constant *F* enabling smoothing and emphasizing of landforms (see Eqs. (5) and (7)). Thus, an adaptation of test results to available independent reference information is possible.

2. Both CSM spatial domain and the explaining relief parameters had to be integrated in the second step. This was done by the application of a hierarchical region-growing segmentation algorithm applied on basic relief parameters considering already existing CSM boundaries. A positive side-effect of segmentation is the reduction of data volume which becomes important for large data sets. Here, the data volume was reduced to one tenth.

3. In fact, clustering is a special case of segmentation in which only neighboring relations in the feature space are considered. Here, the one-dimensional feature spaces of *MBI* or *FPI* were grouped. We have preferred the K-means algorithm because the user can
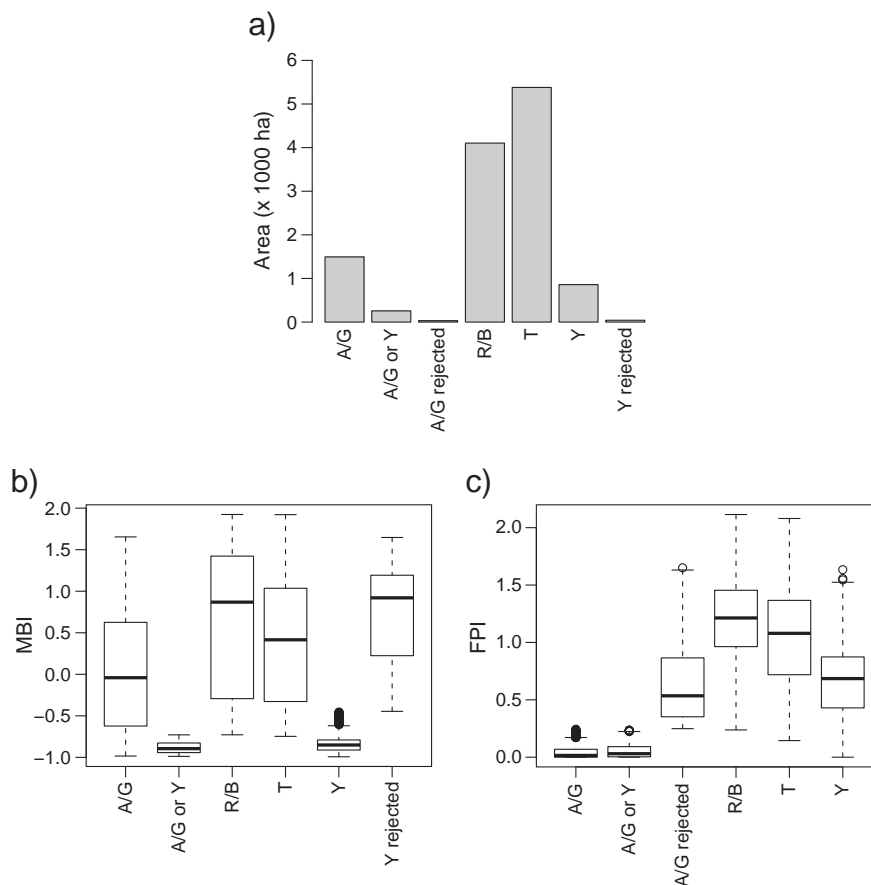


**Fig. 9.** Areas of test result classes (a) as well as corresponding *MBI* (b) and *FPI* box plots (c). A/G and Y consist of confirmed ($C_i$) and positive tested $STO_R$ clusters ($R_{Pi}$; see Fig. 2). T and R/B correspond to negative tested $STO_P$ clusters ($P_i$). The classes 'A/G rejected' and 'Y rejected' are negative tested $STO_R$ clusters ($R_i$) and are equivalent to class ($N_i$). The class 'A/G or Y' represents overlapping distributions of A/G or Y which were tested positive regarding both existence of fluvial and colluvial process domains ($R_{Pi}$).

control the degree of feature space aggregation through the determination of the cluster number. However, the algorithm can be substituted by other approaches which, for instance, enable an automatic determination of cluster number (e.g. Fraley and Raftery, 2007).

4. The applied Kolmogorov Smirnov (KS) test is a traditional goodness-of-fit test calculating the distance $D$ between two empirical cumulative distribution functions (ECDF). Because of its statistical robustness the KS test proved to be suitable for the comparison of soil- and terrain-related distributions. Especially, the insensitiveness to different STO numbers of the comparing distributions is important for the transfer to larger areas. In the future, ECDF-based alternatives to the KS test (see Thas, 2010) will be evaluated.

The presented approach can help to form an impression concerning the terrain-related plausibility of CSM or also legacy soil maps which oftentimes do not contain any accuracy information. Plausibility is expressed here by the distance between reference clusters (RC) and soil unit-related clusters ($P_i$) which should be tested (see Fig. 2). The most crucial point of the workflow is the expert-based RC selection representing "true" process domain of interest. However, in spite of the operation's subjectivity, the selection process itself is done deliberately and controlled and thus, traceable.

In a more general sense, the applied methodology of data integration is suitable for adopting existing data in computer-based systems. The data integration approach is based on the coupling of existing functional hierarchies of thematic geodata (here: a CSM) and multi-scale object structures (Möller et al., 2008). Multi-scale object structures arise from the region-based segmentation of continuous data (here: basic relief parameters). The resulting soil-terrain units enable a supervised disaggregation of heterogeneous soil units (see Wielemaker et al., 2001). Supervised disaggregation means that the applied segmentation algorithm leads to STOs representing a specific geometric aggregation level. The classification of the testing results complies with a semantic disaggregation.

Finally, the test modules may be used for the supervised selection of appropriate samples for automatic classification algorithms (Behrens and Scholten, 2006; Grunwald, 2009; MacMillan, 2008; Scull et al., 2003). Similar to instance selection techniques (Behrens et al., 2008; Schmidt et al., 2008), training data sets might be cleaned from noisy data.

## Acknowledgements

## References

Ad-hoc AG, Boden, 2005. Bodenkundliche Kartieranleitung (KA 5), 5th Edition. E. Schweizerbart'sche Verlagsbuchhandlung, Stuttgart, Germany.

Albrecht, C., Jahn, R., Huwe, B., 2005. Bodensystematik und Bodenklassifikation — Teil 1: Grundbegriffe. Journal of Plant Nutrition Soil Science 168 (1), 7–20.

Baatz, M., Schäpe, A., 2000. Multiresolution segmentation: an optimization approach for high quality multi-scale image segmentation. In: Strobl, J., Blaschke, T. (Eds.), Angewandte Geographische Informationsverarbeitung — Beiträge zum AGIT-Syposium. Vol. 12. Salzburg, pp. 12–23.

Behrens, T., Scholten, T., 2006. A comparison of data-mining techniques in predictive soil mapping. In: Lagacherie, P., A. M., Voltz, M. (Eds.), Digital Soil Mapping — an Introductory Perspective. Vol. 31 of Developments in Soil Science. Elsevier, pp. 353–364.

Behrens, T., Schmidt, K., Scholten, T., 2008. An approach to removing uncertainties in nominal environmental covariates and soil class maps. In: Hartemink, A., McBratney, A., Mendonça Santos, M. (Eds.), Digital Soil Mapping with Limited Data. Springer, pp. 213–224.

Brenning, A., 2008. Statistical geocomputing combining R and SAGA: the example of landslide susceptibility analysis with generalized additive models. In: Böhner, J., Blaschke, T., Montanarella, L. (Eds.), SAGA — seconds out. Vol. 19 of Hamburger Beiträge zur Physischen Geographie und Landschaftsökologie, pp. 23–32.

Bui, E.N., 2004. Soil survey as a knowledge system. Geoderma 120 (1–2), 17–26.

Buol, S., Southard, R., Graham, R., McDaniel, P., 2003. Soil Genesis and Classification, 5th Edition. Iowa State Press, Ames, Iowa.

Davis, J., 2002. Statistics and Data Analysis in Geology. John Wiley & Sons, New York.

Deumlich, D., Thiere, J., Frielinghaus, M., Voelker, L., 1998. MMK characterisation and classification of site conditions in the new federal states of Germany. In: Heineke, H., Eckelmann, W., Thomasson, A., Jones, R., Montanarella, L., Buckley, B. (Eds.), Land Information Systems — Developments for planning the sustainable use of land resources. Vol. 4 of European Soil Bureau Research Report, EUR 17729 EN. Office for Official Publications of the European Communities, Luxembourg, pp. 473–478.

Deumlich, D., Schmidt, R., Sommer, M., 2010. A multiscale soil–landform relationship in the glacial-drift area based on digital terrain analysis and soil attributes. Journal of Plant Nutrition and Soil Science 173 (6), 843–851.

Dobos, E., Hengl, T., 2009. Soil mapping applications. In: Hengl, T., Reuter, H. (Eds.), Geomorphometry — concepts, Software, Applications. Vol. 33 of Developments in Soil Science. Elsevier, pp. 461–479.

Dragut, L., Blaschke, T., 2006. Automated classification of landform elements using object-based image analysis. Geomorphology 81 (3–4), 330–344.

Dragut, L., Eisank, C., 2011. Object representations at multiple scales from digital elevation models. Geomorphology 129 (3–4), 183–189.

Dragut, L., Schauppenlehner, T., Muhar, A., Strobl, J., Blaschke, T., 2009. Optimization of scale and parametrization for terrain segmentation: an application to soil-landscape modeling. Computers and Geosciences 35 (9), 1875–1883.

Fraley, C., Raftery, A., 2007. Model-based methods of classification: using the `mclust` Software in Chemometrics. Journal of Statistical Software 18, 419–430.

Friedrich, K., 1996. Digitale Reliefgliederungsverfahren zur Ableitung bodenkundlich relevanter Flächeneinheiten. Vol. 21 of Frankfurter Geowissenschaftliche Arbeiten. Frankfurt (Main).

Friedrich, K., 1998. Multivariate distance methods for geomorphographic relief classification. In: Heinecke, H., Eckelmann, W., Thomasson, A., Jones, R., Montanarella, L., Buckley, B. (Eds.), Land information systems — developments for planning the sustainable use of land resources. Vol. 4 of European Soil Bureau Research Report, EUR 17729 EN. Office for official publications of the European Communities, Luxembourg, pp. 259–266.

Grunwald, S., 2009. Multi-criteria characterization of recent digital soil mapping and modeling approaches. Geoderma 152 (3/4), 195–207.

Hartmann, K.-J., 2005. Bereitstellung von Informationen der bodenkundlichen Landesaufnahme zur Bewertung von Bodenfunktionen. In: Möller, M., Helbig, H. (Eds.), GIS-gestützte Bewertung von Bodenfunktionen – Datengrundlagen und Lösungsansätze. Wichmann, Heidelberg, pp. 27–34.

Hartmann, K.-J., 2006. Bodenkundliche Basisinformationen. In: Feldhaus, D., Hartmann, K.-J. (Eds.), Bodenbericht 2006 – Böden und Bodeninformationen in Sachsen-Anhalt. Vol. 11 of Mitteilungen zu Geologie und Bergwesen in Sachsen-Anhalt. Landesamt für Geologie und Bergwesen Sachsen-Anhalt, Halle (Saale), pp. 71–88.

Hastie, T., Tibshirani, R., Friedman, J., 2009. The Elements of Statistical Learning: Data Mining, Inference and Prediction, 2nd Edition. Springer Series in Statistics. Springer, New York.

Hengl, T., MacMillan, R., 2009. Geomorphometry – A key to landscape mapping and modelling. In: Hengl, T., Reuter, H. (Eds.), Geomorphometry - Concepts, Software, Applications. Vol. 33 of Developments in Soil Science. Elsevier, pp. 433–460.

Hutchinson, M., 1989. A new procedure for gridding elevation and stream line data with automatic removal of spurious pits. Journal of Hydrology 106, 211–232.

ISO 19138, 2006. Geographic information: Data quality measures. Tech. rep. International Organization for Standardization, Geneve, Switzerland.

IUSS Working Group WRB, 2006. World Reference Base for Soil Resources, 2nd ed. FAO, Rome.

MacMillan, R., 2008. Experiences with applied DSM: Protocol, availability, quality and capacity building. In: Hartemink, A., McBratney, A., Mendonça Santos, M. (Eds.), Digital Soil Mapping with limited data. Springer, pp. 113–135.

MacMillan, R., Shary, P., 2009. Landforms and landform elements in geomorphometry. In: Hengl, T., Reuter, H. (Eds.), Geomorphometry - Concepts, Software, Applications. Vol. 33 of Developments in Soil Science. Elsevier, pp. 227–254.

Minár, J., Evans, I., 2008. Elementary forms for land surface segmentation: The theoretical basis of terrain analysis and geomorphological mapping. Geomorphology 95 (3–4), 236–259.

Möller, M., 2008. Scale-specific derivation of thematic basic data for landscape analysis. Ph.D. thesis. University of Tübingen, Tübingen, Germany, in German.

Möller, M., Volk, M., Friedrich, K., Lymburner, L., 2008. Placing soil genesis and transport processes into a landscape context: A multi-scale terrain analysis approach. Journal of Plant Nutrition and Soil Science 171, 419–430.

Müller, E., Volk, M., 2001. History of landscape assessment. In: Krönert, R., Steinhardt, U., Volk, M. (Eds.), Landscape balance and landscape assessment. Springer, Berlin, pp. 23–46.

Olaya, V., Conrad, O., 2009. Geomorphometry in SAGA. In: Hengl, T., Reuter, H. (Eds.), Geomorphometry - Concepts, Software, Applications. Vol. 33 of Developments in Soil Science. Elsevier, pp. 293–308.

Quinn, P., Beven, K., Lamb, R., 1995. The ln(a/tan b) index: How to calculate it and how to use ist within the TOPMODEL framework. Hydrological Processes 9, 161–182.

Schmidt, K., Behrens, T., Scholten, T., 2008. Instance selection and classification tree analysis for large spatial datasets in digital soil mapping. Geoderma 146 (1–2), 138–146.

Scull, P., Franklin, J., Chadwick, O.A., McArthur, D., 2003. Predictive soil mapping: A review. Progress in Physical Geography 27 (2), 171–197.

Thas, O., 2010. Comparing Distributions. Springer Series in Statistics, Springer, New York.

Wielemaker, W.G., de Bruin, S., Epema, G.F., Veldkamp, A., 2001. Significance and application of the multi-hierarchical landsystem in soil mapping. Catena 43 (1), 15–34.

Zevenbergen, L., Thorne, C., 1987. Quantitaive analysis of land surface topography. Earth Surface Processes and Landforms 12, 12–56.