

# Uncertainty propagation in VNIR reflectance spectroscopy soil organic carbon mapping

L. Brodský\*, R. Vašát, A. Klement, T. Zádorová, O. Jakšík

Department of Soil Science and Soil Protection, Faculty of Agrobiological Sciences, Czech University of Life Sciences Prague, Kamýcká 129, Prague 6, Czech Republic

## ARTICLE INFO

### Article history:

Received 23 November 2011  
Received in revised form 7 November 2012  
Accepted 19 November 2012  
Available online 20 December 2012

### Keywords:

Soil organic carbon  
VNIR spectroscopy  
Partial least squares regression  
Digital soil mapping  
Uncertainty

## ABSTRACT

Visible and near infrared (VNIR) diffuse reflectance spectroscopy (DRS) offers high potential as a fast and accurate proximal soil sensing technique for soil carbon estimation. The objective of this study is to evaluate the use of VNIR soil spectroscopy for mapping soil organic carbon (SOC) spatial distribution on a 100 ha arable field strongly affected by erosion. The analysis was performed in two main steps: firstly, we focused on the uncertainty in the VNIR spectroscopy regression model (PLSR) under varying number and locations of training samples from which an optimal number of input samples were selected; secondly, we analysed uncertainty propagation in the coupled PLSR and spatial prediction for the selected optimal number of training samples.

The PLSR quality parameters are changing exponentially with increasing number of input training samples. The PLSR model constructed using only 37 samples provided a good predictive capability with  $R^2$  over 0.7 and RPD over 1.5. The uncertainty of the final map as expressed by mean standard deviation, lowered 3 times when the number of input training samples changed from 37 to 128. The accuracy of the soil map was assessed through the uncertainty propagation analysis for the purpose of evaluating how the uncertainties were partially associated with the PLSR model and partially coming from the spatial prediction propagate to the final output map. We conclude that the PLSR predictions caused lower uncertainty in comparison with uncertainty coming from spatial predictions by kriging algorithm. The study confirms that the SOC prediction as made by VNIR spectral characteristics is a powerful tool, which together with digital soil mapping techniques (DSM) provides the basis for high resolution field scale mapping.

© 2012 Elsevier B.V. All rights reserved.

## 1. Introduction

Soils originating from loess are highly susceptible to erosion and re-deposition due to their physical characteristics. A soil mosaic comprising of various degradation and accumulation phases of original Chernozem (soil unit defined by the World Reference Base; FAO, 2006) is common on cultivated land as a result of early deforestation and intensive agricultural use (Terhorst, 2000). The low density (e.g. soil organic carbon) and fine particles (clay and silt) are removed preferentially from the profile due to a selective process of soil erosion (Lal, 2005). Mapping spatial patterns of soils in landscapes affected by erosion focuses mostly on topographic analysis of the widely available digital elevation model (Florinsky et al., 2002; McKenzie and Ryan, 1999; Odeh et al., 1995; Zádorová et al., 2011) and conventional laboratory analysis.

Visible and near infrared (VNIR) diffuse reflectance spectroscopy (DRS) offers high potential as a fast and accurate proximal soil sensing technique for principal soil property estimation (Stenberg and Rossel, 2010; Stenberg et al., 2010). VNIR spectroscopy is being considered to complement or even replace in some cases conventional soil analysis. The research, development and application of the technique in soil

science are clearly documented by an increasing number of studies (e.g. Ben-Dor and Banin, 1995; Dalal and Henry, 1985; Dematte, 2004; Shepherd and Walsh, 2004; Stenberg et al., 2010; Viscarra Rossel et al., 2006). Dalal and Henry (1985) emphasized the technique as rapid, non-destructive and providing simultaneous measurements of a number of soil properties. Viscarra Rossel et al. (2006), Brown (2007) and others regard the technique as inexpensive. The commonly reported drawback of the technique is that the VNIR predictions do not achieve the laboratory precision of chemical analysis (e.g. Brown, 2007). On the other hand, Viscarra Rossel et al. (2006) claimed that DRS spectroscopy can sometimes be more accurate than the conventional analysis. Validation of the technique against laboratory chemical analysis shows varying results. Stenberg et al., 2010 provide a tabulated overview of the validation results from different studies. The most frequently estimated property is SOC. Factors negatively influencing the precision of the VNIR predictions are discussed by Stenberg et al. (2010) and were identified as follows: varying soil texture – sand content, soil structure, moisture, mineralogy, and others. It should also be noted that the different studies applied different reference laboratory analysis procedures. Although, the predictive capacity of VNIR spectroscopy is limited due to a number of factors negatively influencing the estimation prediction, which require more attention, there is high potential of DRS due to the efficiency of the information acquisition. Viscarra Rossel and McBratney (2008) proposed DRS as a tool for digital

\* Corresponding author. Tel.: +420 224382632; fax: +420 234381836.  
E-mail address: [brodsky@af.czu.cz](mailto:brodsky@af.czu.cz) (L. Brodský).

soil mapping (DSM). They reviewed the costs and benefits of soil DRS for DSM and identified two initial investments, a spectrometer and the development of a spectral library.

As the individual steps of VNIR diffuse reflectance spectroscopy are being continually developed, there is a need for quantifying and evaluating the associated uncertainties in the context of the application examples. Much attention is paid to the multivariate calibrations of the reflectance spectra and consequently to the soil property prediction. A large number of comparisons of different algorithms can be found in the scientific literature. Mouazen et al. (2010) provided a comparison of principal component regression (PCR), partial least squares regression (PLSR), and back propagation neural network (BPNN). Viscarra Rossel and Behrens (2010) compared multiple linear regression (MLR), PLSR, multivariate adaptive regression splines (MARS), support vector machine (SVM), random forest (RF) boosted trees (BT), and artificial neural networks (ANN).

Although the artificial intelligence group of models (ANN, SVM, RF) usually outperformed the others, the PLSR technique provides the most competitive predictions and is successfully applied in many studies for the prediction of various soil properties (Bartholomeus et al., 2011; Janik et al., 2007; Mouazen et al., 2010; Reeves, 2010; Stenberg, 2010; Terhoeven-Urselmans et al., 2010; Viscarra Rossel et al., 2006). The advantages of PLSR are that it is relatively computationally fast, models are more interpretable, and the method is not as sensitive to over-fitting as ANN or SVM. Also the PLSR copes with data containing a large number of collinear predictor variables and the model explains more variance in the response with fewer components (Mouazen et al., 2010; Stenberg et al., 2010; Viscarra Rossel and Behrens, 2010). PLSR takes the soil information within the internal computation process and makes use of the correlation between spectra and the soil data, unlike PCR which computes a set of eigenvectors independently and then performs regression. Viscarra Rossel et al. (2006) report that “by fitting a PLSR model, one hopes to find a few PLSR factors that explain most of the variation in both predictors and responses”. During the calibration process it is important to determine the optimal number of factors to retain in the calibration model; using too few factors can leave important parts of structure the un-modelled, and using too many factors draws too much measurement noise into the calibration model (Björsvik and Martens, 2008).

Most of the VNIR reflectance spectroscopy studies in soil science were processed under laboratory conditions, while some other studies prepare for so-called mobile (on-line or on-the-go) digital soil mapping (Adamchuk et al., 2011; Mouazen et al., 2007) in field conditions without the necessity to manually collect samples. Mobile mapping also provides very high density of spatial sampling.

Once the soil property predictions by VNIR spectroscopy are made, spatial predictions by various kriging algorithms can be applied to draw the maps, following the digital soil mapping concept (McBratney et al., 2003). A large number of interpolation methods have been developed and elaborated. In general, these can be: (a) interpolation relying solely on point observations of the target variable; and (b) interpolation based on regression of the target variable on spatially exhaustive auxiliary information.

The whole procedure of combined VNIR spectroscopy and spatial predictions is fairly straightforward and can save costs on laboratory analyses. The resulting soil property maps will be inaccurate given uncertainties that arise in different steps within the mapping procedure such as soil sampling, spectra collecting, building the VNIR prediction model (number of training samples and model parameters), and final spatial prediction. The accuracy of the resulting soil maps can be insufficient for the intended purpose, e.g. field-level optimized management, therefore it is important to be able to quantify the accuracy of the maps.

The objective of this study is to test the ability of VNIR soil spectroscopy for SOC mapping on the steep sloping arable lands strongly affected by water erosion. We attempt here to perform an uncertainty

propagation analysis of SOC mapping with the use of VNIR spectroscopy coupled with spatial predictions and to quantify the uncertainty associated with different steps within the mapping procedure (i.e. the uncertainty rising from PLSR model and predictions, and from ordinary kriging predictions).

The analysis was performed in two main steps: firstly, we focused on the uncertainty in the VNIR spectroscopy regression model (PLSR) under varying number and locations of training samples from which an optimal number of input samples were selected; secondly, we analysed uncertainty propagation in the coupled PLSR and spatial prediction for the selected optimal number of training samples.

## 2. Materials and methods

### 2.1. Soil sampling and laboratory analysis

The study was conducted on an arable field with an area of 100 ha located in South-eastern Czech Republic (South Moravia, Brumovice municipality) representing a rough terrain landscape (Fig. 1a). The digital elevation model (DEM) is plotted with 2 m spacing contour lines derived from a topographic map. The target site represents an intensively exploited piece of arable land with commonly cultivated crops. According to WRB (FAO, 2006) the main soil units are defined as Haplic Chernozem and Regosol. A detailed description of the study area can be found in Zádorová et al. (2011).

The basic protocol of VNIR spectroscopy application is assembled as follows:

- (i) soil samples collection,
- (ii) spectra collection,
- (iii) laboratory analysis of training samples,
- (iv) statistical modelling of the relationship between soil properties and VNIR spectra,
- (v) statistical prediction at locations for which only spectra information is available.

Soil samples were extensively collected from the top soil to a depth of 20 cm, to simulate the situation which is close to mobile mapping (Fig. 1b). A total of 781 soil samples were collected in a way which covers the geographic distribution of the area the most uniformly. The sampling strategy can be described as design based systematic sampling on a regular grid with varying density, where higher density is applied in areas of strong erosion and deposition impact, and lower density on the upper flat plain which is not affected so much by erosion. Initially the North-West deposition area was sampled with higher density to investigate the magnitude of spatial variation of SOC, after which the rest of the 100 ha field was sampled uniformly. In the next step the soil samples were split into three groups, (training, testing and prediction data) in order to get data for model calibration (training data), independent validation of the resulting model (testing data) and prediction of SOC (prediction data). The sub-setting into the three groups with portion of 152 (training set), 50 (testing set), and 579 (prediction set) was done by judgement sampling when soil samples were split by eye to get the most uniform spatial distribution within each group. The training and testing soil samples were undertaken for laboratory analysis of SOC and spectra measurement, while the prediction samples underwent spectra measurement only. The size of the training and the test set was determined to be high enough to allow stable regression modelling as well as validation, while the training set was decided to be in the magnitude of about 150 to allow also for reduction evaluating the influence of the size of training set on PLSR model performance.

The samples were air-dried, ground and sieved to 2 mm. Oxidisable organic carbon content was determined oxidimetrically by a modified Tyurin method (Pospíšil, 1964).

Soil spectra were collected by standard procedure under laboratory conditions with a FieldSpec 3 (350–2500 nm) instrument equipped

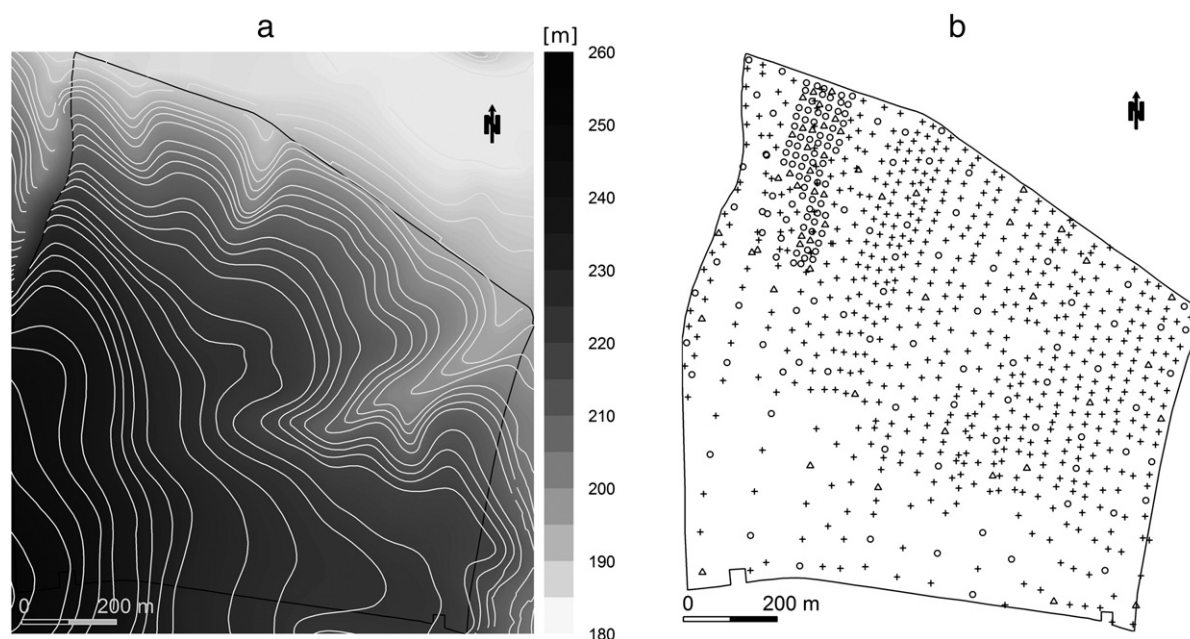


Fig. 1. Brumovice field digital elevation model with contours (a) and sampling sites allocations: ○ training samples, Δ test samples and + prediction samples (b).

with a High Intensity Contact Probe. The spectrometer was calibrated by a Spectralon standard white reference panel (Halon). The calibration was checked after every 10th sample. Each soil sample was mixed before spectra collection to avoid preference distribution of particles.

All computations, graphs and figures were made with R software tool for statistical computing (R Development Core Team, 2011).

Summary statistics of the soil laboratory analysis (Table 1), as well as VNIR spectra, were evaluated for representatives of the two data sets. Fig. 3 presents histograms of the training (a) and testing (b) data set.

## 2.2. PLSR modelling

The pls package: Principal Component and Partial Least Square Regression in R (Mevik and Wehrens, 2007), using *pls* function, was employed to model the relationship between SOC and the spectra. Original spectral data from the instrument were converted into spectral reflectance and smoothed by cubic spline prior to PLSR analysis. Several other data manipulations were tested, transformed to optical density units by  $\log_{10}(\text{Reflectance})$  and 1st derivative, but positive effects on the PLSR modelling and predictions were negligible.

Partial least-squares regression (PLSR) initially with a full training set (152) and maximum 20 components was used to calibrate the spectral data with the reference laboratory soil data. The model was

independently validated against the testing set by Root Mean Square Error of Predictions (RMSEP). The RMSEP was used as criteria to select “the best” calibration model. The procedure is primarily to select a model (number of factors) with lowest RMSEP. In this study we also considered relative change of RMSEP with number of factors compared to minimum RMSEP. Viscarra Rossel et al. (2006) applied a statistical approach with *t*-test and defined alpha. Beside RMSEP, explained variance in relation to number of factors was evaluated. The other quality assessment parameters were coefficient of multiple determinations  $R^2$  and residual prediction deviation (RPD). The first four PLSR factor loading weight factors, considered to be indicators of the correlation between the spectra and soil constituents, were plotted for the qualitative interpretation. Model performance was also compared to published results. PLSR modelling in the initial evaluation was also tested with a reduced training set (100, 50 and 20), by random sub-setting without replacement, to see the performance change as preparation for the automatic batch processing in the next simulations.

## 2.3. Geostatistical modelling

Geostatistical procedures, ordinary kriging (OK) predictions as well as spatial stochastic simulations with OK, were performed with gstat package (Pebesma, 2004). For R scripting we used further packages: *maptools* (Lewin-Koh and Bivand, 2011) and *sp* (Bivand et al., 2008; Pebesma and Bivand, 2005).

Semi-variogram modelling was processed by the classic method of moment variogram estimate. After initial plotting, an automated fitting procedure *fit.variogram* using weighted least squares was applied (weighted by number of points and distance). The variogram model was fitted sequentially with an exponential, spherical and liner model from which the exponential model had the lowest SSErr (sum of squares for error) and hence it was kept since it describes the spatial variation of SOC the most adequately. The automatic fitted variogram was then modified by eye if necessary.

The number of nearest observations that are used for prediction (neighbourhood size) was set as 30 for both OK prediction and conditioned stochastic simulation with OK. The stochastic simulation was conditioned to the SOC observations.

Table 1  
Summary statistics of SOC.

Data set	Mean (%)	Median (%)	Range (%)	Standard deviation (%)	Coefficient of variation (%)	Skewness
Full (202 samples)	1.11	1.06	0.31–1.92	0.34	30.8	0.22
PLSR-training (152 samples)	1.15	1.13	0.41–1.92	0.36	30.9	−0.01
PLSR-testing (50 samples)	1.10	1.04	0.31–1.81	0.35	31.7	0.17

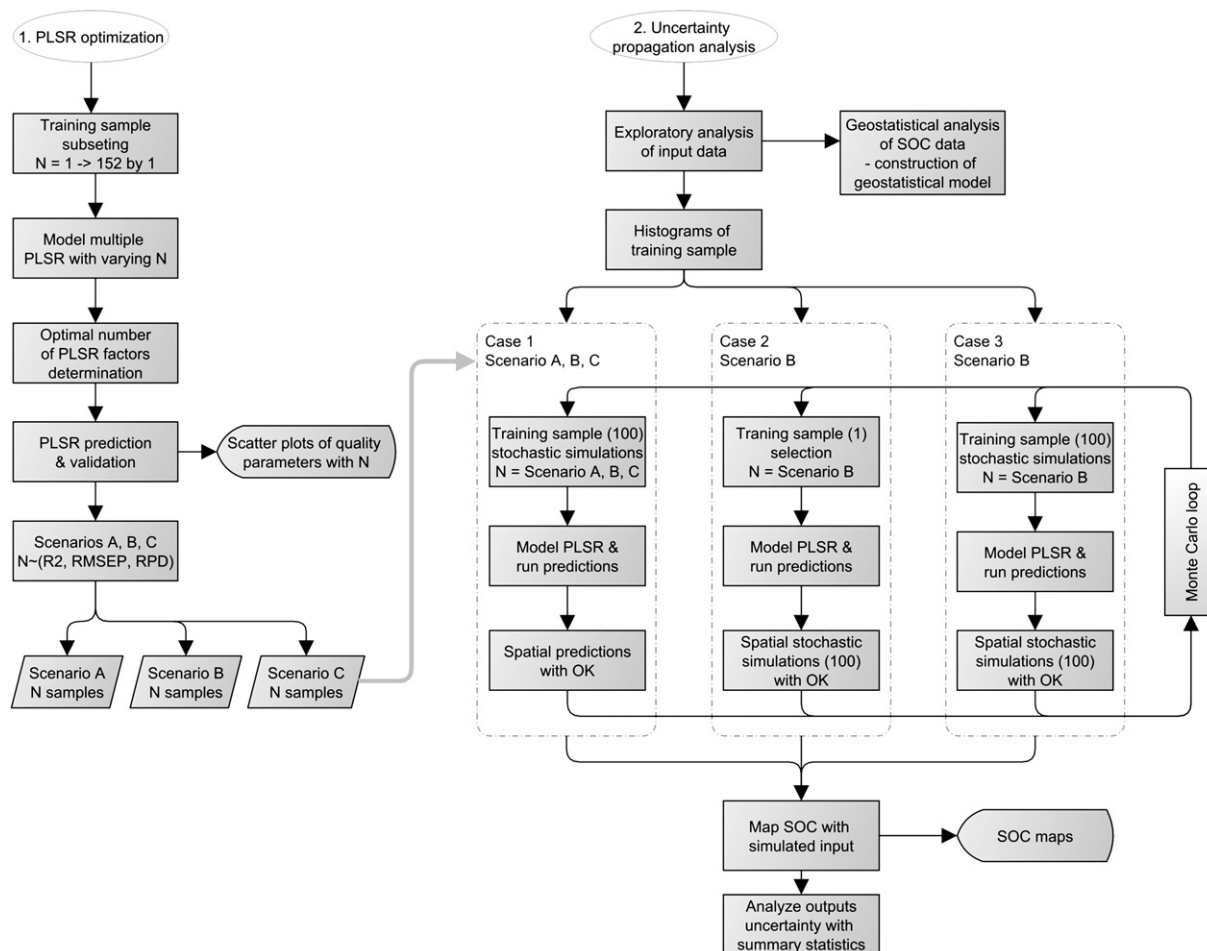


Fig. 2. Schema of the uncertainty analysis workflow.

## 2.4. Uncertainty analysis

The analysis firstly focused on the uncertainty in the VNIR spectroscopy regression model under varying number and locations of training samples and subsequently on the selection of the optimal number of input samples. Secondly we aimed at the uncertainty propagation in the coupled PLSR and spatial predictions procedure under the selected optimal number of training samples, Fig. 2.

### 2.4.1. PLSR uncertainty under varying number of training samples

Firstly, uncertainty in the statistical inference model was assessed by a varying number of training samples in order to optimize the PLSR model. Construction of multiple PLSR models was run with a changing number of input training samples  $N$  from 11 to 152 by 1. The number of factors in PLSR model was limited according to initial PLSR analysis to the maximum plausible number. The best model, defined by the number of factors, in one run was then selected by

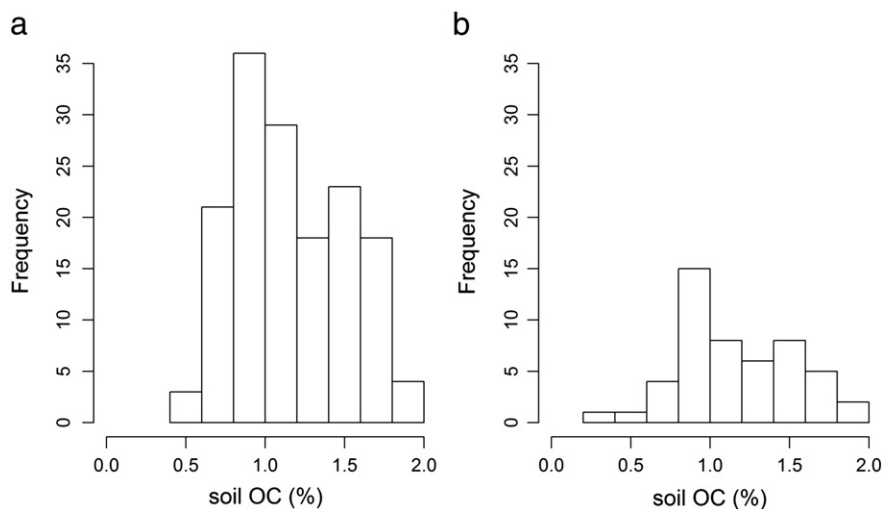


Fig. 3. Histogram of training (a) and testing (b) data set of SOC.



the RMSEP parameter as calculated from independent test data. Coefficient of multiple determination  $R^2$  and Residual Prediction Deviation (RPD) statistics were also calculated to evaluate the model performance. RPD is calculated as the ratio of the standard deviation of reference (training data set) soil property values to the RMSE of prediction (Viscarra Rossel, 2007). Next, three different scenarios, different number and allocation of training samples, were selected according to the calculated and determined model accuracy parameters. The general criterion in the analysis was to find model with the determined accuracy and minimum training samples required. The number of training samples is an important parameter that directly relates to labour of the laboratory analysis. The three scenarios are defined to provide trade-off between model quality and cost.

Scenario A: moderate quality PLSR model with  $R^2 \geq 0.6$ ,  $RPD \geq 1.5$ ;  
 Scenario B: good quality PLSR model with  $R^2 \geq 0.7$  and  $RPD \geq 1.5$ ;  
 Scenario C: the best quality PLSR model from the given data set with  $R^2 \geq 0.8$  and  $RPD \geq 2.0$ .

#### 2.4.2. Uncertainty propagation of the SOC map

Soil map accuracy can be evaluated by collecting a large number of independent validation data with negligible measurement error and comparing these with the predictions from the soil map. Another approach can be to utilize the uncertainty propagation analysis as described comprehensively by Heuvelink (1998). In our study, the issue of how uncertainties in VNIR spectroscopy modelling and the spatial interpolation will propagate to the output map is presented. Since uncertainty analysis is a very complex task, the Monte Carlo method (e.g. Bishop et al., 2006; Heuvelink et al., 2007) was applied. The Monte Carlo method builds a loop over the coupled PLSR and spatial predictions and can be briefly described in a few steps:

1. repeat  $i$  times ( $i \geq 100$ ):
  - i. generate a realisation from the probability distribution of the uncertain inputs using a pseudo-random number generator;
  - ii. run the coupled models (PLSR-OK) with the inputs and store the result;
2. compute summary statistics (mean and standard deviation) of the  $N$  model outputs.

Initially, mapping of SOC from available laboratory analysis (without help of VNIR spectroscopy) was performed by OK to later compare the uncertainty with a full data set (including VNIR spectroscopy data). Next, the accuracy assessment of the resulting maps through an uncertainty analysis was made by three separate calculations which are illustrated in Fig. 2.

Secondly, the uncertainty propagation from the PLSR modelling and predictions was analysed for the three scenarios defined with a different number of training samples (A, B and C). The multiple training data set (100 subsets) for each scenario was constructed by sample selection of  $N$  members using a pseudo-random number generator. Sampling was constrained to the original sample data set to reproduce the statistical distribution. Each of the 100 simulations was followed by the PLSR modelling with SOC prediction and spatial prediction by OK. The final uncertainty was consequently calculated as a standard deviation map for each scenario from the 100 runs.

Thirdly, spatial stochastic simulation (100 realisations) of SOC using OK algorithm was run on one PLSR predicted data using the  $N$  number of samples defined according to B scenario to assess the uncertainty associated with spatial prediction.

Finally, a combined simulation of training sample from scenario B and stochastic spatial simulation (with OK) were run resulting in  $100 \times 100$  combinations in the Monte Carlo loop. As an indicator of accuracy of final SOC maps we used the standard deviation computed from all maps.

### 3. Results

#### 3.1. Soil organic carbon spatial distribution

The soil chemical laboratory analysis (202 samples) proved an expected wide variation of SOC in the erosion field. The summary statistics of the full, PLSR-training and PLSR-testing data set are presented in Table 1.

The semi-variogram model constructed using all available laboratory SOC data, indicates moderately strong spatial auto-correlation of SOC with nugget = 0.03, sill = 0.09, and range = 150 m for the exponential model, fitted by eye. The proportion of nugget to sill is high (33%). The pattern of SOC map (Fig. 4a) is highly smoothed. A more detailed mapping is indicated in the NW part of the area where the sampling density is higher. This coincides with the NW valley affected by erosion and deposition. Uncertainty of the prediction map (Fig. 4b) is presented as average kriging standard deviation in Table 2. Unfortunately the level of uncertainty is high – nearly at the level of SOC variation in the field (Table 1).

#### 3.2. Initial PLSR analysis

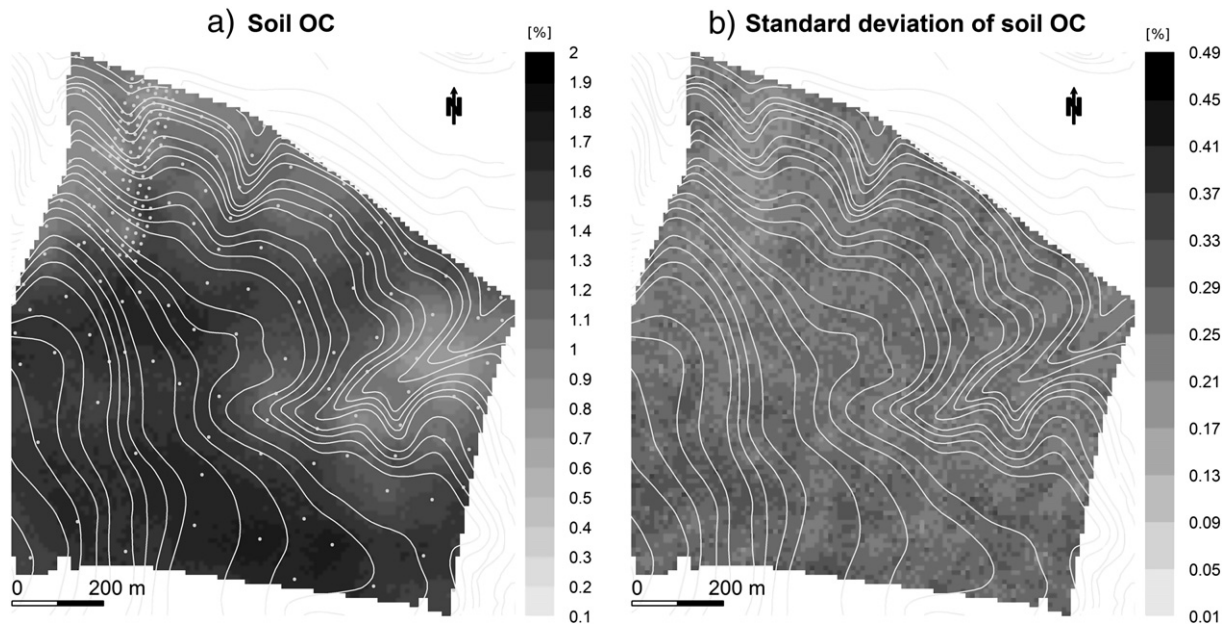
Soil spectra measured on 781 samples are very similar with little variation and minor absorption features. Here we focus on evaluation of the initial PLSR analysis of SOC from the full training data set (202 samples). Fig. 5 shows cross-validated root mean square error of prediction against the number of factors and the first four PLSR factor loading weight vectors. The first plot shows a rapid decrease of RMSEP when the factor equals to three and more. The global minimum of RMSEP was found to be 0.137 and belongs to the 12 factors run, however, the difference in RMSEP for 10 factors run was negligibly small (2.9% of the minimum). RMSEP increases again when the number of factors equals to 12 and more. We consider that factors higher than 10 keep mostly noise. Furthermore, 99.9% of variation in spectral training data is explained when the number of factors equals to four and 100% for ten factors.

The loading weight vectors from the PLSR decomposition (Fig. 5b) indicate positive and negative peaks, which occur at specific wavelengths. The first factors show a negative loading value and indicate inference, while the other three factors indicate positive peaks – considered as correlation between spectra and the soil components of interest. The validation statistics for PLSR model with ten factors is as follows:  $RMSEP = 0.141\%$ ,  $R^2 = 0.85$  and  $RPD = 2.4$ . The result shows high correlation of the soil spectra and laboratory analysis.

#### 3.3. PLSR uncertainty analysis under varying number of samples

In this part of the study we attempt to evaluate the performance of PLSR model under a changing number of training samples. Uncertainty of the model was assessed by three statistical parameters:  $R^2$ , RMSEP and RPD. The results from the independent validation of the PLSR models against the changing number of training samples are presented in Fig. 6 (○ 10-factor PLSR model). The RMSEP is exponentially decreasing and stabilizes around 0.17%, while  $R^2$  and RPD are increasing exponentially with the number of input training samples. Please note that the plots show dispersed trends rather than narrow lines. This effect is attributed to the use of pseudo-random generator for the training sample selection in the analysis.

The results revealed that the PLSR model constructed from only 20 training soil samples provided moderate predictive capability as defined in scenario A. A good quality PLSR model as defined in scenario B was achieved using 37 training samples, while the model with  $R^2$  over 0.8 and RPD over 2.0 was achieved taking 128 training samples (scenario C). The selection of the appropriate number of training samples ( $N$ ) from the graphs is made as subjective due to the high point dispersion; hence the three scenarios were evaluated by taking the



**Fig. 4.** Mean (a) and standard deviation (b) maps of SOC (%) from 202 soil laboratory analysis without the use of VNIR spectroscopy data.

lowest number of training samples and corresponding model quality criteria. The summary statistics of SOC for the three scenarios training subsets is provided in Table 2.

Viscarra Rossel et al. (2006) found that “a four-factor model caused an insignificant decrease in accuracy” and suggest that a more parsimonious model than the one corresponding with the lowest RMSE may be appropriate. An additional test with 4-factor PLSR model was run (Fig. 6, points plotted as +). As a result of this test the 4-factor model showed similar quality parameters for scenario A. The three statistical parameters ( $R^2$ ; RMSEP and RPD) reached saturation from about 60 training samples with the quality indicators worse than the 10-factor model.

### 3.4. Uncertainty propagation to the soil map

Next, scenarios A, B and C with  $N = 20, 37$  and  $128$  were analysed to calculate uncertainty within the mapping process. Simulations of different samples selection from the training data set were run in the 100 iterations. As a result there are 100 maps for the uncertainty analysis. These were summarized by mean and standard deviation, where the mean can be interpreted as the best estimate and the standard deviation as a measure of the accuracy of the estimate. The summary kriging map of SOC with uncertainty for  $N = 37$  is presented in Fig. 7. The mean maps for  $N = 20, 37$  and  $128$  showed very similar patterns with a mean value  $1.33\%$  and a negligibly varying range of SOC. The spatial patterns of the standard deviation maps were also

highly similar but with varying values of standard deviation. The mean standard deviation of the simulations with  $N = 20$  was  $0.135\%$ , with  $N = 37$  it was  $0.092\%$  and with  $N = 128$  it was only  $0.028\%$  (Table 3). Areas with higher prediction uncertainty coincide mostly with locations where the content of SOC is low.

The semi-variogram model constructed from the exhaustive data (781 samples) indicates strong spatial autocorrelation of SOC. The semi-variogram parameters were as follows: nugget =  $0.02$ , sill =  $0.12$ , and range =  $100$  m for the exponential model, fitted by eye. The nugget to sill ratio equals to  $16\%$ , and lowered considerably compared to the initial semi-variogram modelling with 202 samples.

To compare the results of the previous simulations, 100 ordinary kriging simulations of the SOC were run, where the realisations were conditioned to the SOC predictions by the best PLSR model.

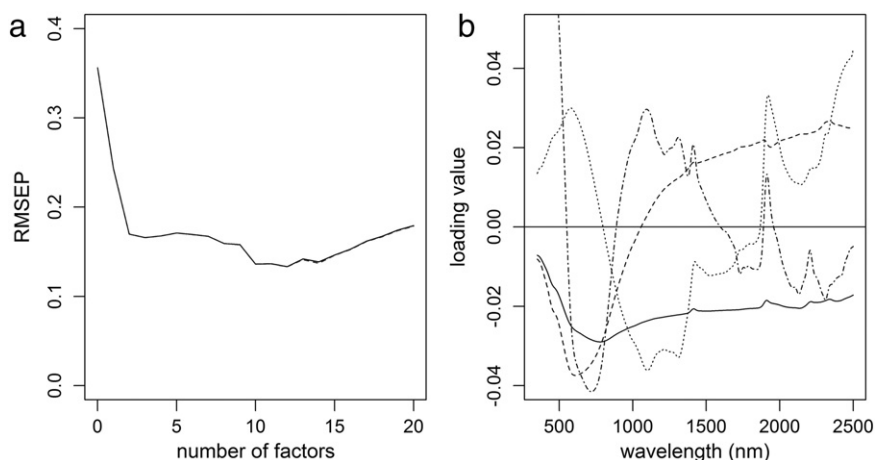
The kriging map, calculated as a mean of the 100 realisations, showed a similar spatial pattern as presented in Fig. 7a. The mean value of the kriging simulations was  $1.34\%$ . The difference between the two mean maps is small. The spatial pattern of the standard deviation map (Fig. 8) partly reveals the pattern of the soil sampling design, while the lower sampling density in the SW part resulted in higher prediction uncertainty. The mean value of standard deviation in the map is  $0.149\%$ , which is considerably higher compared to the value coming from the PLSR simulations ( $0.092\%$ ). The uncertainty in the spatial prediction is higher than uncertainty of PLSR prediction under varying sample locations. Moreover, the difference was statistically confirmed by  $t$ -test ( $t = 54.67$  and  $p\text{-value} = 2.2e - 16$ ). The mean value of standard deviation is also slightly higher even than scenario A ( $0.135\%$ ). This conveys that the uncertainty in the spatial prediction is higher than the uncertainty coming from PLSR with a reduced number of training samples.

Finally, we combined simulations of training sample selection for PLSR (model scenario B) and stochastic spatial simulations based on OK (Fig. 2). These simulations resulted in  $100 \times 100$  combinations. The output of uncertainty analysis is 10,000 maps from which mean and standard deviation are presented in Fig. 9. The mean map can be interpreted as the best estimate and the standard deviation as a measure of the accuracy of the estimate. It is clear that the spatial pattern of the realisations reproduces those as shown in Fig. 7a. The differences in the maps are small. The kriging prediction map in Fig. 7 is somewhat smoother than the map in Fig. 9. The standard deviation of the estimate was  $0.16\%$ . The map pattern is similar to that of

**Table 2**

Summary statistics of SOC for scenarios A, B and C (training subsets).

Data set	Mean (%)	Median (%)	Range (%)	Standard deviation (%)	Coefficient of variation (%)	Skewness
Scenario A (20 samples)	1.10	1.10	0.50–1.64	0.29	26.6	0.03
Scenario B (37 samples)	1.12	1.16	0.50–1.71	0.32	27.1	−0.01
Scenario C (128 samples)	1.16	1.09	0.41–1.92	0.36	30.6	0.23



**Fig. 5.** PLSR cross-validated root mean square error of prediction (RMSEP) against the number of factors (a) and the first four PLSR factor loading weight vectors: 1 – solid, 2 – dashed, 3 – dotted, and 4 – dot-dash lines (b).

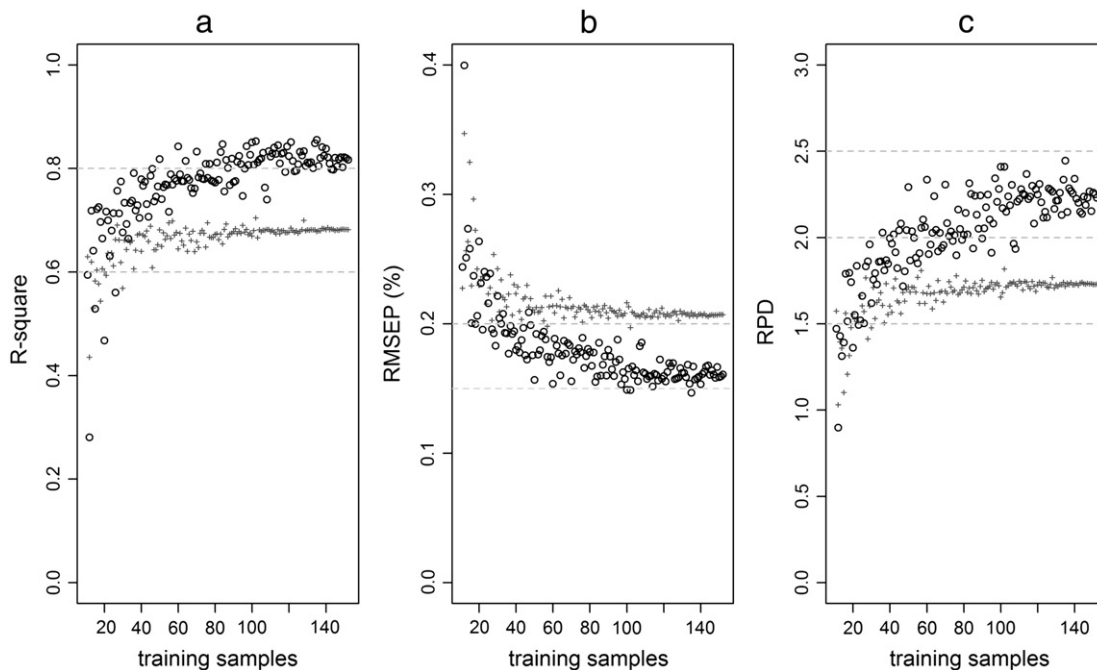
the standard deviation map from OK spatial predictions (Fig. 9a), while there are locally higher values that propagated from PLSR predictions (Fig. 7b).

#### 4. Discussion

The map of SOC spatial distribution (Fig. 7) shows areas affected by erosion and related processes to different extents. Bright parts indicate a low content of SOC, which correspond to the back-slope segments with high slope and convex curvature areas covered mainly by Regosols and eroded Chernozems. The SW flat upper part of the plot shows a high SOC content which corresponds with low erosion activity. The concave side valleys and the toe slope are characterized by 0.6–1.5% SOC content. These higher values occur as a result of the accumulation of both organic and mineral material from the neighbouring slopes. The SOC map indicates an actual state of SOC spatial distribution, which is invaluable information in soil erosion assessment.

By comparing the maps of SOC (Fig. 4) generated from only available samples with laboratory chemical analysis and with the use of VNIR spectroscopy (Fig. 7), it is evident that the latter shows higher spatial variation. The maps were qualitatively evaluated by a soil erosion expert having local knowledge of the field with the conclusion that the former map provides insufficient detail. Sampling density certainly needs to be higher than 202 samples for a 100 ha area with local spatial variation of this complexity. Moreover, the 202 samples for laboratory chemical analysis is too high an investment to map one arable field. This implicates the need of VNIR spectroscopy as an inexpensive method to provide soil information for high resolution soil mapping. This finding supports the discussed potential of DRS and its application to digital soil mapping in Viscarra Rossel and McBratney (2008).

The other question related to the cost of laboratory chemical analysis is how many training samples are needed for VNIR spectroscopy. The results, shown in Fig. 6, demonstrate that it is a trade-off between quality of the prediction model and the cost. Another option could be to use a large spectral library and subset it, but the quality then



**Fig. 6.** Scatter plots of PLSR parameters  $R^2$  (a), RMSEP (b) and RPD (c) plotted against number of training samples from 141 simulations (○ 10-factor model, + 4-factor model).



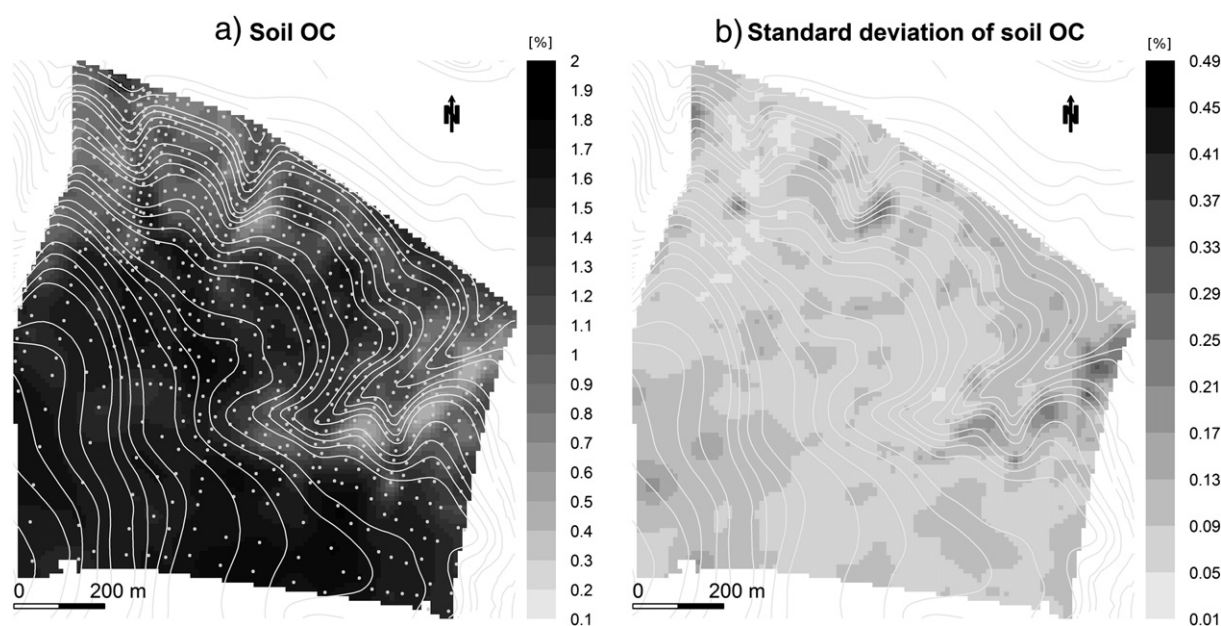


Fig. 7. Mean (a) and standard deviation (b) maps of SOC (%) from 100 PLSR sample simulations (scenario B).

strongly depends on representativeness of the training samples. Here we consider only local samples collected from one geologically homogeneous field.

According to the scenarios A, B and C, moderate quality of VNIR prediction model can be constructed with as little as 20 training samples collected from the field. A good quality model, as defined in the materials and methods, was constructed with only 37 training samples. However, sample selection also plays an important role as can be seen from the dispersion of the points in Fig. 6. This uncertainty is revealed by the standard deviation map (Fig. 7b) where higher values appear in areas where the SOC is low, for which fewer samples were available in the training data set. The number of required samples for the suggested model quality parameters may change in other fields or other environmental conditions.

Uncertainty propagation analysis reveals the main source of uncertainty in the process of SOC mapping. The maps of kriging standard deviation together with mean value in Table 3 document the results that the highest uncertainty is present in the maps created from samples with only laboratory chemical analysis, without the use of VNIR spectroscopy. Simulated sample selection for PLSR modelling generally caused a lower uncertainty in the maps than spatial predictions by kriging. A possibility to decrease the uncertainty from spatial prediction would be to increase spatial sampling density by for

instance using a mobile mapping technique (Adamchuk et al., 2011; Mouazen et al., 2007). Nevertheless, uncertainty in the soil spectra from soil moisture would have to be tackled. The impact of soil moisture on reflectance spectra is non-linear and can be higher than differences in reflectance due to varying SOC, hence can destroy the SOC analysis. Another possibility would be to apply spatial sampling optimization, e.g. spatial coverage sample method or spatial simulated annealing (Vašát et al., 2010). The goal of the latter approach is to minimize the sample size while keeping the average kriging variance values of the tested soil variables below a given threshold. Combined PLSR and spatial prediction simulations showed a slightly higher value of uncertainty, mean standard deviation, than the separate PLSR and spatial simulations with OK (Table 3). This implies that the uncertainties from PLSR and kriging in the whole mapping

**Table 3**  
Analysis of uncertainty – summary statistics.

Uncertainty analysis	Number of PLSR training samples	Number of samples for kriging	Mean kriging standard deviation (%)
OK (without VNIR predictions)	–	202	0.253
PLSR predictions (Scenario A)	20	781	0.135
PLSR predictions (Scenario B)	37	781	0.092
PLSR predictions (Scenario C)	128	781	0.028
Spatial predictions OK (Scenario B)	37	781	0.149
Combined PLSR and OK simulations (Scenario B)	37	781	0.160

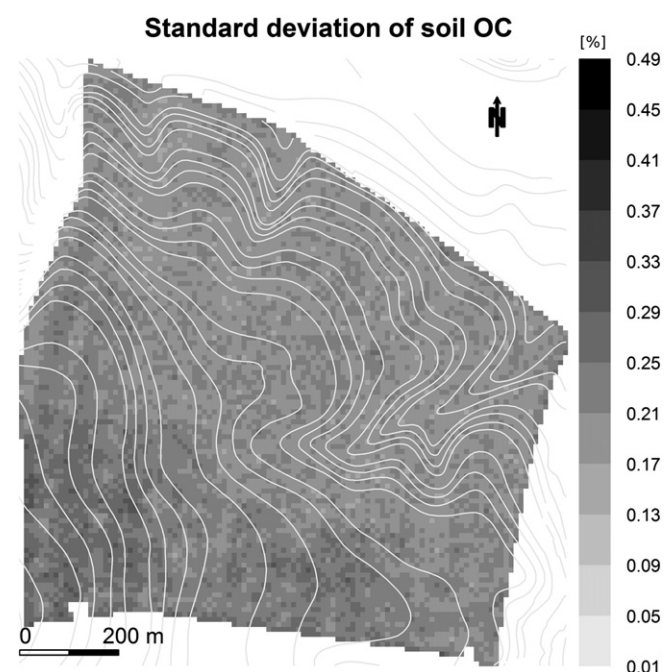


Fig. 8. Ordinary kriging standard deviation map of SOC from 100 spatial simulations (scenario B).



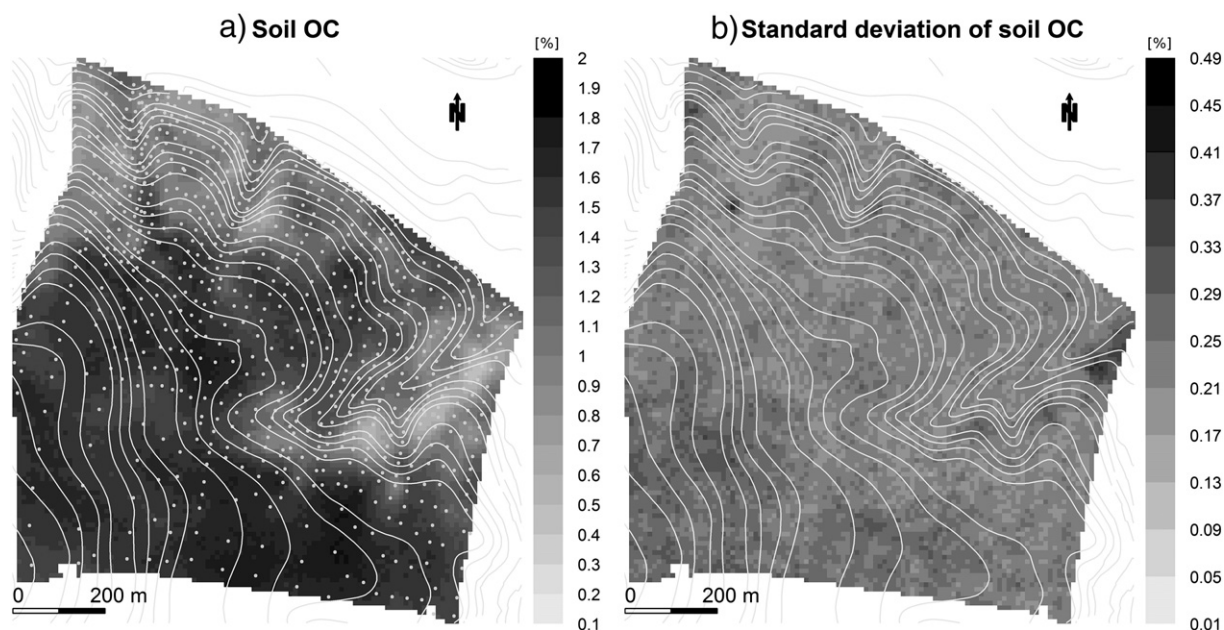


Fig. 9. Mean (a) and standard deviation (b) maps of SOC from 100×100 combined PLSR and spatial simulations (scenario B).

procedure have additive effects on the final map. The locally higher standard deviation caused by PLSR modelling, higher than that coming from OK, propagated to the final result (Fig. 9b). To reduce this effect of higher uncertainty of PLSR, which coincide mostly with locations where the SOC content is low, one should use an improved training set also covering the lower quartile well.

## 5. Conclusions

The study confirms that the SOC prediction as made by VNIR spectral characteristics is an effective tool, which together with digital soil mapping techniques provides the basis for high resolution field scale mapping. In the case of a geologically homogenous area, a strong correlation between soil spectra and training laboratory data was found.

The PLSR model quality indicators are changing exponentially with increasing number of training samples. The relative stability of the indicators appears to be achieved using approximately 80 training samples. However, the model constructed from only 37 training soil samples provided results with  $R^2$  over 0.7 and RPD over 1.5. A model with  $R^2$  over 0.8 and RPD over 2.0 was achieved using 128 training samples. The uncertainty of the final map as expressed by the mean standard deviation is lowering 3 times when the number of input training samples is changing from 37 to 128. Previous findings document the trade-off between VNIR spectroscopy prediction model quality and cost.

Accuracy of the resulting SOC map was evaluated through the uncertainty assessment analysis to evaluate how the uncertainties partially associated with the PLSR model and partially coming from the spatial prediction will propagate to the final output. The uncertainty associated with PLSR modelling resulted in lower uncertainty in the predicted maps compared to the uncertainty coming from spatial predictions as made by kriging algorithm. This conveys a very good predictive capability of the VNIR spectroscopy. Uncertainties from both, PLSR and kriging, in the whole mapping procedure had additive effects on the final map.

## Acknowledgements

The authors acknowledge the financial support of the Czech Science Foundation grant No. 526/09/1762, 526/08/0434 and by the Ministry of Education, Youth and Sports of the Czech Republic, Project No. MSM 6046070901.

## References

- Adamchuk, V.I., Viscarra Rossel, R.A., Marx, D.B., Samal, A.K., 2011. Using targeted sampling to process multivariate soil sensing data. *Geoderma* 163 (1–2), 63–73.
- Bartholomeus, H., Kooistra, L., Stevens, A., Van Leeuwen, M., Van Wesemael, B., Ben-Dor, E., Tychon, B., 2011. SOC mapping of partially vegetated agricultural fields with imaging spectroscopy. *International Journal of Applied Earth Observation and Geoinformation* 13 (1), 81–88.
- Ben-Dor, E., Banin, A., 1995. Near-infrared analysis as rapid method to simultaneously evaluate several soil properties. *Soil Science Society of America Journal* 59, 364–372.
- Bishop, T.F.A., Minasny, B., McBratney, A.B., 2006. Uncertainty analysis for soil-terrain models. *International Journal of Geographical Information Science* 20, 117–134.
- Bivand, R.S., Pebesma, E.J., Gomez-Rubio, V., 2008. *Applied Spatial Data Analysis with R*. Springer, NY.
- Björsvik, H.R., Martens, H., 2008. *Data Analysis: Calibration of NIR Instruments PLS Regression*. Handbook of Near-Infrared Analysis, 3rd edition, CRC Press, Taylor & Francis Group, Boca Raton, FL, 189–205.
- Brown, D.J., 2007. Using a global VNIR soil-spectral library for local characterization and landscape modelling in a 2nd-order Uganda watershed. *Geoderma* 140, 444–453.
- Dalal, R.C., Henry, R.J., 1985. Simultaneous determination of moisture, OC, and total nitrogen by near infrared reflectance spectrophotometry. *Soil Science Society of America Journal* 50 (1), 16–19.
- Dematte, J., 2004. Visible–NIR reflectance: a new approach on soil evaluation. *Geoderma* 121 (1–2), 95–112.
- FAO, 2006. *World Reference Base for Soil Resources*. FAO, Rome.
- Florinsky, I.V., Eilers, R.G., Manning, G.R., Fuller, L.G., 2002. Prediction of soil properties by digital terrain modelling. *Environmental Modelling and Software* 17, 295–311.
- Heuvelink, G.B.M., 1998. *Error Propagation in Environmental Modelling with GIS*. Taylor & Francis, London. (127 pp.).
- Heuvelink, G.B.M., Brown, J.D., Van Loon, E.E., 2007. A probabilistic framework for representing and simulating uncertain environmental variables. *International Journal of Geographical Information Science* 21 (5), 497–513.
- Janik, L.J., Merry, R.H., Forrester, S.T., Lanyon, D.M., Rawson, A., 2007. Rapid prediction of soil water retention using mid infrared spectroscopy. *Soil Science Society of America Journal* 71 (2), 507.
- Lal, R., 2005. Forest soils and carbon sequestration. *Forest Ecology and Management* 220, 242–258.
- Lewin-Koh, N.J., Bivand, R., 2011. Maptools: tools for reading and handling spatial objects, R package version 0.8–10. <http://CRAN.R-project.org/package=maptools>.
- McBratney, A.B., Mendonça Santos, M.L., Minasny, B., 2003. On digital soil mapping. *Geoderma* 117 (1–2), 3–52.
- McKenzie, N.J., Ryan, P.J., 1999. Spatial prediction of soil properties using environmental correlation. *Geoderma* 89, 67–94.
- Mevik, B.-H., Wehrens, R., 2007. The pls package: principal component and partial least squares regression in R. *Journal of Statistical Software* 18 (2), 1–24.
- Mouazen, A.M., Maleki, M.R., De Baerdemaeker, J., Ramon, H., 2007. On-line measurement of some selected soil properties using a VIS–NIR sensor. *Soil and Tillage Research* 93 (1), 13–27.
- Mouazen, A.M., Kuang, B., De Baerdemaeker, J., Ramon, H., 2010. Comparison among principal component, partial least squares and back propagation neural network analyses for accuracy of measurement of selected soil properties with visible and near infrared spectroscopy. *Geoderma* 158 (1–2), 23–31.

- Odeh, I.O.A., McBratney, A.B., Chittleborough, D.J., 1995. Further results on prediction of soil properties from terrain attributes: heterotopic cokriging and regression-kriging. *Geoderma* 67, 215–226.
- Pebesma, E.J., 2004. Multivariable geostatistics in S: the gstat package. *Computers and Geosciences* 30, 683–691.
- Pebesma, E.J., Bivand, R.S., 2005. Classes and methods for spatial data in R. *R News* 5 (2) (<http://cran.r-project.org/doc/Rnews/>).
- Pospíšil, F., 1964. Fractionation of humus substances of several soil types in Czechoslovakia. *Rostlinná Výroba* 10, 567–579.
- R Development Core Team, 2011. R: a language and environment for statistical computing. R Foundation for Statistical Computing. (<http://www.R-project.org>).
- Reeves, J.B., 2010. Near- versus mid-infrared diffuse reflectance spectroscopy for soil analysis emphasizing carbon and laboratory versus on-site analysis: where are we and what needs to be done? *Geoderma* 158 (1–2), 3–14.
- Shepherd, K.D., Walsh, M.G., 2004. Diffuse reflectance spectroscopy for rapid soil analysis. In: Lal, Rattan (Ed.), *Encyclopedia of Soil Science*.
- Stenberg, B., 2010. Effects of soil sample pretreatments and standardised rewetting as interacted with sand classes on Vis–NIR predictions of clay and SOC. *Geoderma* 158 (1–2), 15–22.
- Stenberg, B., Viscarra Rossel, R.A., 2010. Diffuse reflectance spectroscopy for high-resolution soil sensing. In: Viscarra Rossel, R.A., McBratney, A.B., Minasny, B. (Eds.), *Media*, vol. 1. Springer, Netherlands, pp. 29–47.
- Stenberg, B., Viscarra Rossel, R.A., Mouazen, A.M., Wetterlind, J., 2010. Visible and near infrared spectroscopy in soil science. *Advances in Agronomy* 107, 163–215.
- Terhoeven-Urselmans, T., Vagen, T.-G., Spaargaren, O., Shepherd, K.D., 2010. Prediction of soil fertility properties from a globally distributed soil mid-infrared spectral library. *Soil Science Society of America Journal* 74 (5), 1792–1799.
- Terhorst, B., 2000. The influence of Pleistocene landforms on soil-forming processes and soil distribution in a loess landscape of Baden-Württemberg. *Catena* 41, 165–179.
- Vašát, R., Heuvelink, G.B.M., Borůvka, L., 2010. Sampling design optimization for multi-variate soil mapping. *Geoderma* 155 (3–4), 147–153.
- Viscarra Rossel, R.A., 2007. Robust modelling of soil diffuse reflectance spectra by "bagging-partial least squares regression". *Journal of Near Infrared Spectroscopy* 15 (2), 39–47.
- Viscarra Rossel, R.A., Behrens, T., 2010. Using data mining to model and interpret soil diffuse reflectance spectra. *Geoderma* 158 (1–2), 46–54.
- Viscarra Rossel, R.A., McBratney, A.B., 2008. Diffuse reflectance spectroscopy as a tool for digital soil mapping. In: Hartemink, Alfred E., McBratney, Alex B., Mendonça-Santos, Maria de Lourdes (Eds.), *Digital Soil Mapping with Limited Data. Part III*. Springer, pp. 165–172.
- Viscarra Rossel, R.A., Walvoort, D.J.J., McBratney, A.B., Janik, L.J., Skjemstad, J.O., 2006. Visible, near infrared, mid infrared or combined diffuse reflectance spectroscopy for simultaneous assessment of various soil properties. *Geoderma* 131, 59–75.
- Zádorová, T., Penížek, V., Šefrna, L., Rohošková, M., Borůvka, L., 2011. Spatial delineation of OC-rich Colluvial soils in Chernozem regions by Terrain analysis and fuzzy classification. *Catena* 85, 22–33.