

Digital soil mapping using data with different accuracy levels



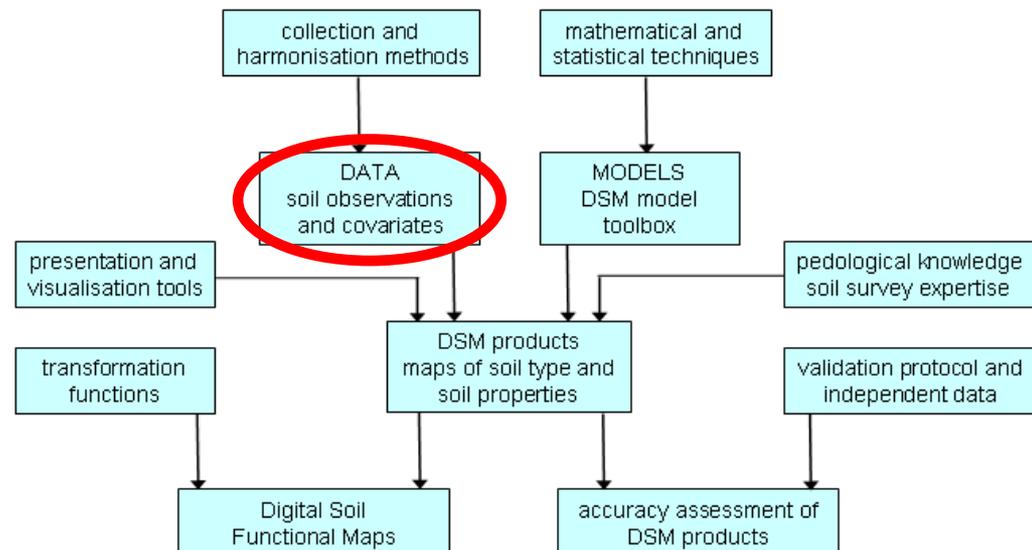
World Soil Information

Gerard Heuvelink, Dick Brus, Tom Hengl,
Bas Kempen, Johan Leenaars and Maria
Ruiperez-Gonzalez

Soil profile data are key to DSM

- They are used to **calibrate** soil prediction models that predict soil properties from covariates
- They are used to **condition** predictions to nearby observations using kriging

Digital Soil Mapper



But soil profile data are not without error

- Lab data can have substantial **measurement errors**
- There are **quality differences** between labs
- Field data ('educated guesses') tend to be **less accurate** than lab data
- Many soil 'observations' are **measured indirectly** (such as through soil spectroscopy, e.g. PLSR)
- Some of us make use of **pseudo-observations**
- In future the use of **volunteered soil information** (crowd-sourcing) will grow, but these data may not be very reliable



World Soil Information

SOILINFO APP

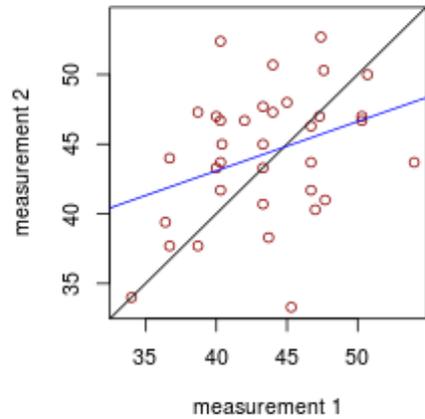
Providing free access to soil data across borders

<http://soilinfo.isric.org>

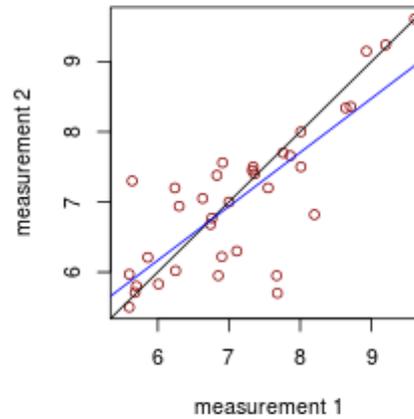


Quality of lab data: look and shudder

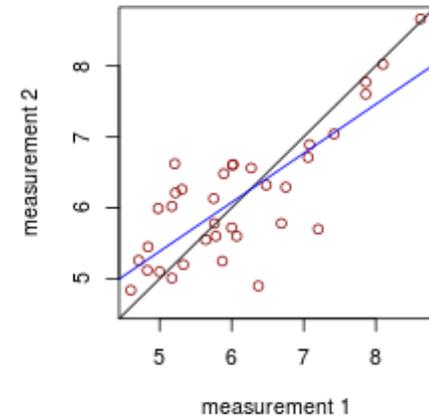
Sat_w.x



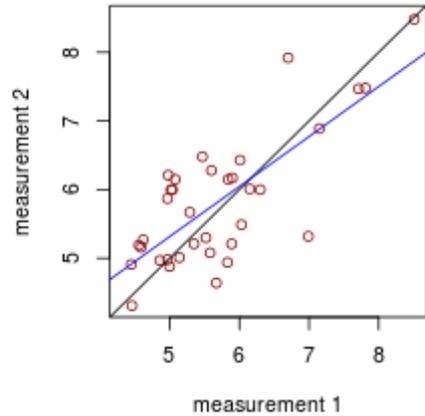
pHw.x



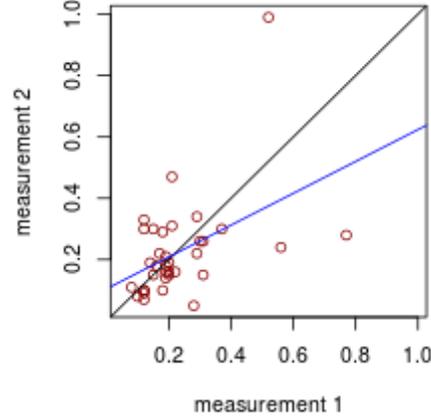
pHCaCl.x



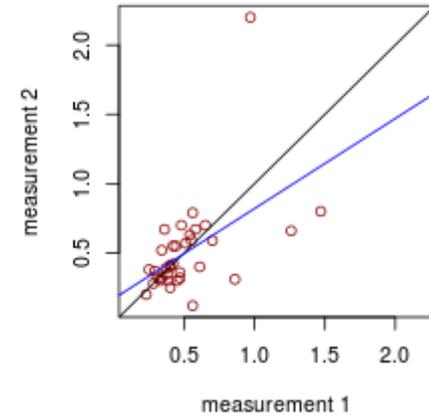
pHKCl.x



CE_susp.x



CE_sat.x



The geostatistical approach can account for measurement errors (KED case)

$$Z(s) = m(s) + \varepsilon(s) = \sum_{j=0}^p \beta_j \cdot x_j(s) + \varepsilon(s)$$

$$Y(s_i) = Z(s_i) + \delta(s_i)$$

measurement errors with mean vector μ
and variance-covariance matrix V

$$\hat{Z}(s_0) = E[Z(s_0) | Y(s_i) = y(s_i), i = 1 \cdots n]$$



Result that goes back to Delhomme (1978)

$$\hat{\beta} = (X^T (C + V)^{-1} X)^{-1} X^T (C + V)^{-1} (y - \mu)$$

$$\text{Var}(\hat{\beta} - \beta) = (X^T (C + V)^{-1} X)^{-1}$$

$$\hat{Z}(s_0) = (x_0 + X(X^T (C + V)^{-1} X)^{-1} (x_0 - X^T (C + V)^{-1} c_0))^T \cdot (C + V)^{-1} (y - \mu)$$

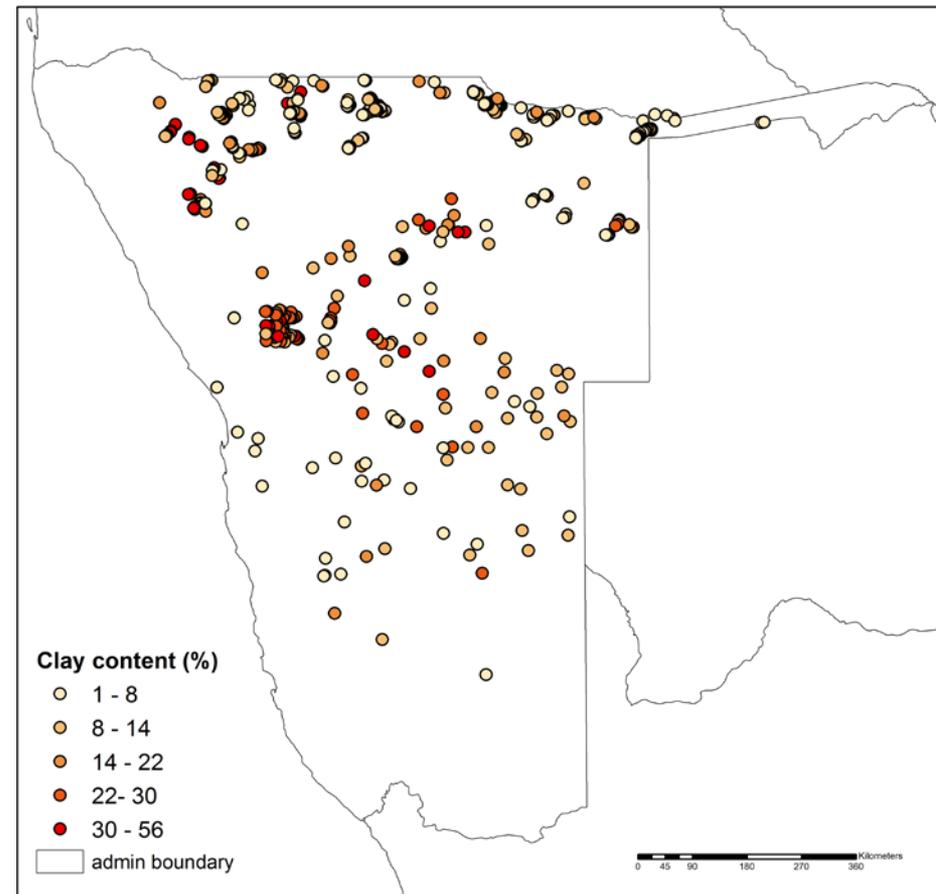
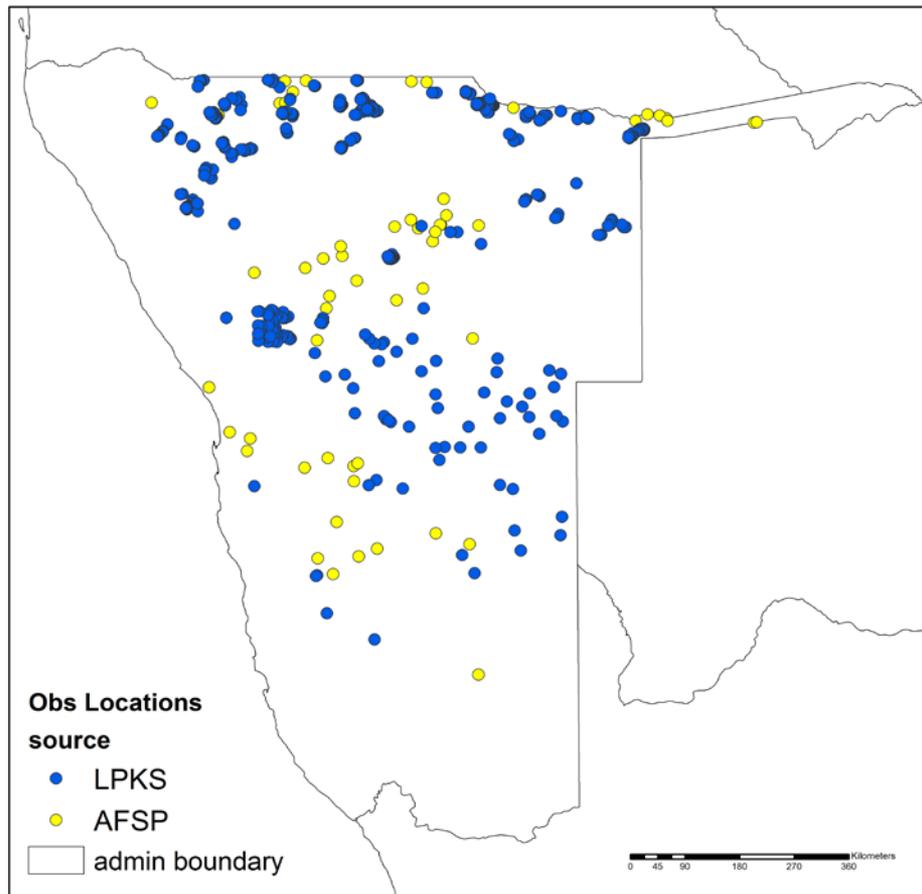
$$\begin{aligned} \text{Var}(\hat{Z}(s_0) - Z(s_0)) &= c_{00} - c_0^T \cdot (C + V)^{-1} \cdot c_0 \\ &+ (x_0 - X^T (C + V)^{-1} c_0)^T (X^T (C + V)^{-1} X)^{-1} \\ &\cdot (x_0 - X^T (C + V)^{-1} c_0) \end{aligned}$$



Example: mapping topsoil clay content for Namibia

LPKS = LandPKS database, field estimates of soil texture (by texture class)

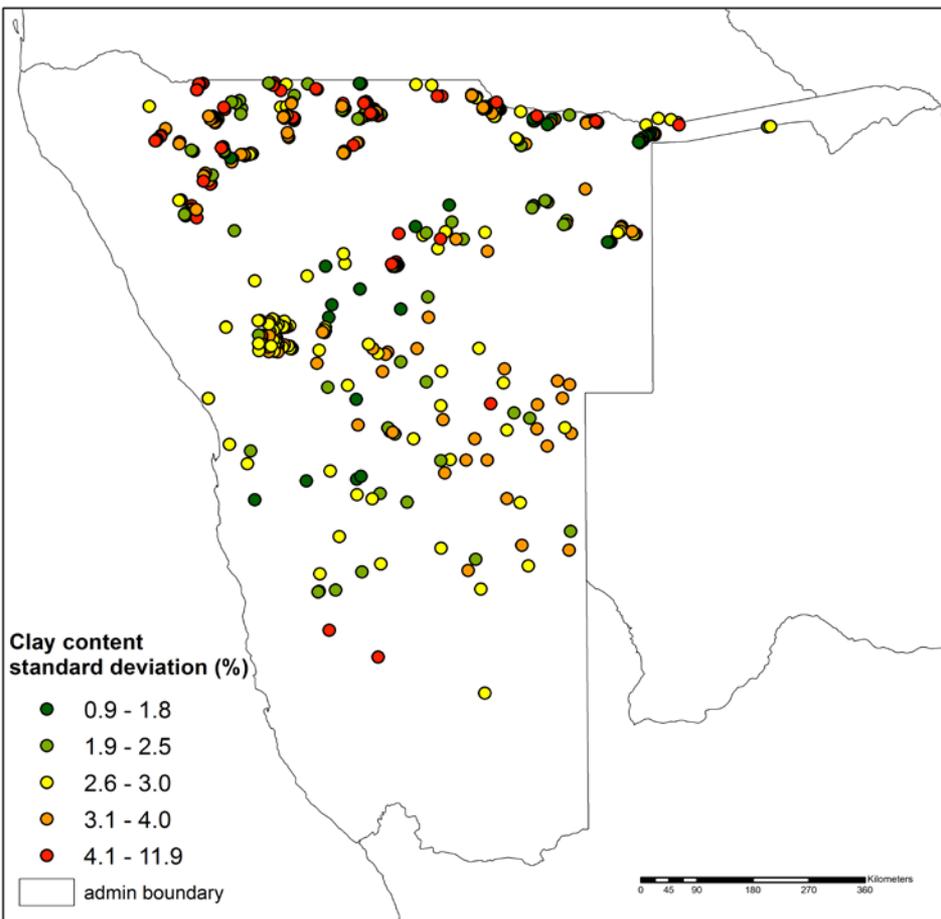
AFSP = Africa Soil Profiles database (merge of numerous legacy soil datasets)



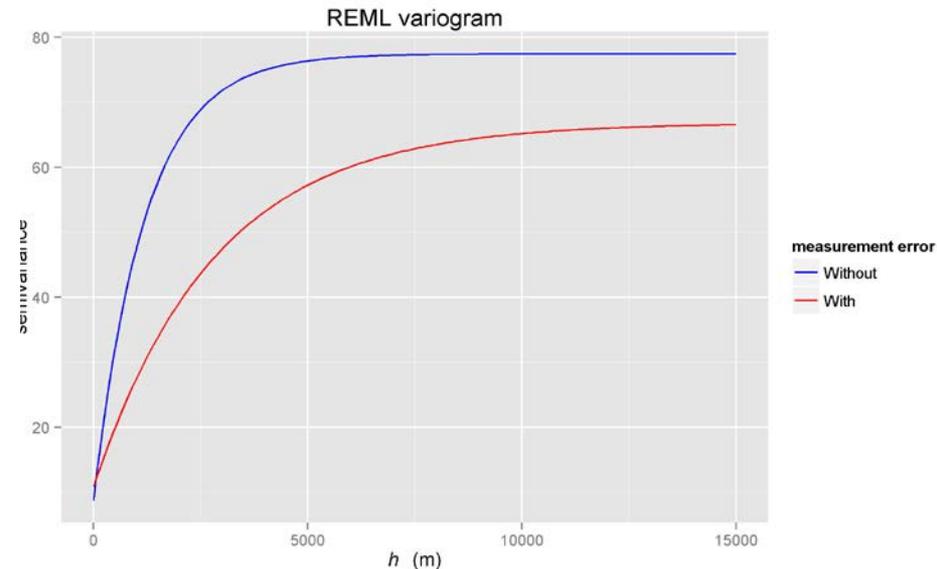
AFSP data **four accuracy levels** (depending on source credibility)

LPKS data accuracy based on **variability within soil texture class** (GSIF TT2tri function in R)

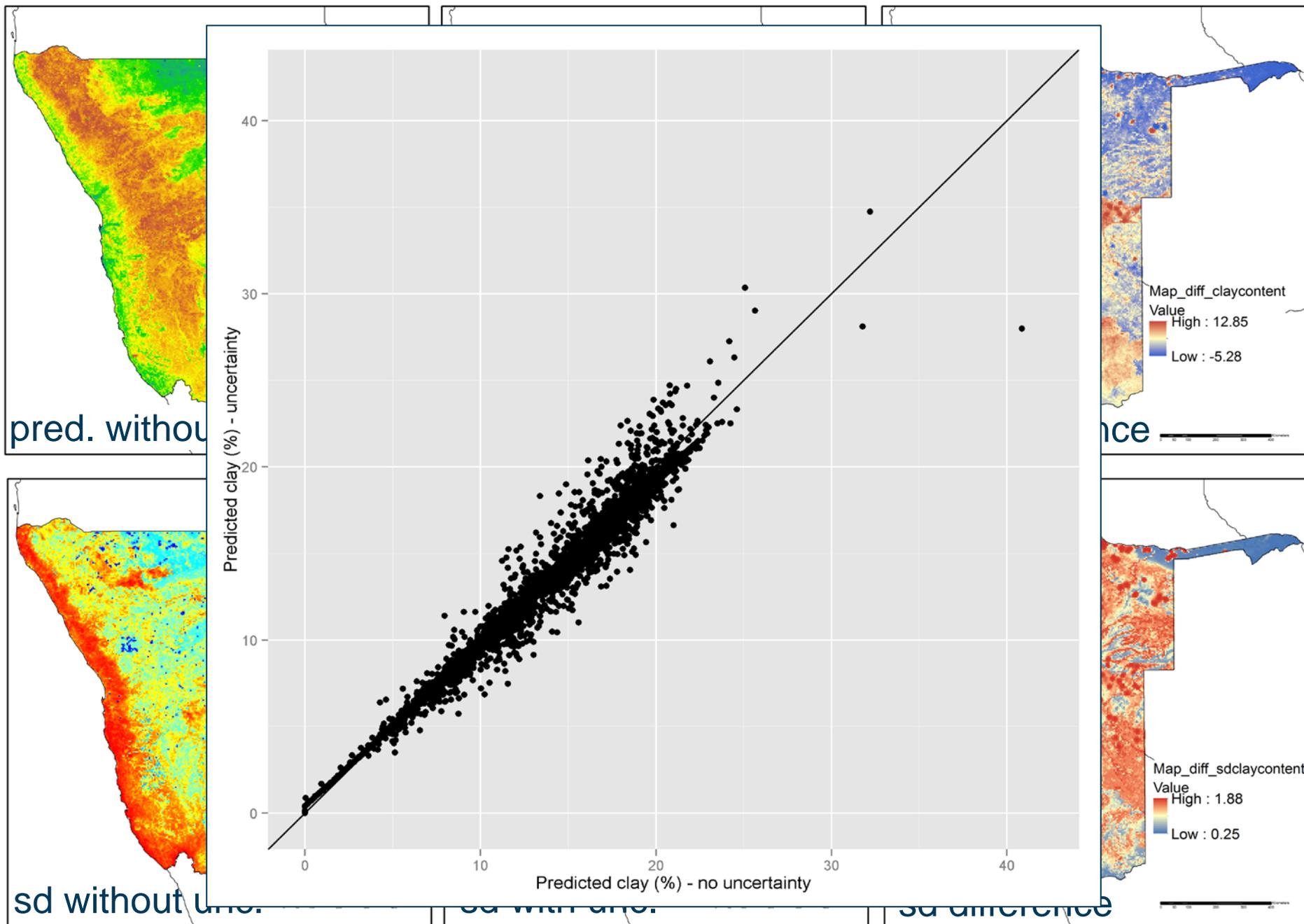
clay sd ranges from 0.9% to 11.9%



	Beta_GLS (without uncertainty)	Beta_GLS (including uncertainty)
Intercept	8.596	8.835
ASSDAC3	-0.063	-0.057
VBFMRG5	-0.0075	-0.0073
T03MSD3	0.227	0.208



Resulting maps: meaningful differences



This was just an illustrative example, possible **extensions**:

- Include (known) **systematic differences** between sub-datasets
- Include **unknown** systematic differences by representing these as random errors that are **perfectly correlated** within a sub-dataset
- Include **serial correlation** between errors (instrument drift, anchoring effect)
- Take different soil data accuracy levels into account for **variogram estimation** (including its uncertainty, using a Markov chain Monte Carlo approach)
- Estimate part of the **measurement error parameters** (i.e. elements of μ and V) from the data



Also important and challenging:

- How to include differences in soil data accuracies in **machine-learning** algorithms?
- One possibility is to assign **weights**, but how large should these weights be?



NIH Public Access

Author Manuscript

Stat Anal Data Min. Author manuscript; available in PMC 2014 February 03.

Published in final edited form as:

Stat Anal Data Min. 2013 December 1; 6(6): 496–505. doi:10.1002/sam.11196.

A Weighted Random Forests Approach to Improve Predictive Performance

Stacey J Winham^{1,§}, Robert R Freimuth¹, and Joanna M Biernacka^{1,2,§}

¹Department of Health Sciences Research, Mayo Clinic, 200 First Street Southwest, Rochester, MN 55905 USA

Concluding remarks

- Soil measurements are **not error-free**
- Measurement error **can be taken into account** in DSM, so why don't we do it?
- Perhaps it is because often we do not know how **accurate** the soil measurements are?
- But this is **not true** for data such as derived from soil spectroscopy, and why don't we routinely send **replicates** to the laboratory (without telling the lab)?
- **We can do so much better**. And we should. Is there anyone in this room who has not wasted valuable time on trying to fit models to 'poor' (rubbish) data?



Submission of
pre-conference
workshop
proposals
**1 November,
2016**

Abstract
submission
deadline
**1 February,
2017**

Early-bird
conference
registration
1 April, 2017

Submission of
special issue
manuscripts
1 October, 2017

Evolution Proximal Soil Sensing Soil Monitoring

Please join us in Wageningen, The Netherlands!

26 June – 1 July, 2017

Come share and learn about the latest developments in the field and join the festive 25th anniversary of pedometrics.

Stay up to date at
www.pedometrics2017.org

SoilCares 



ΠΕΔΟ
METRICS

