

Soil Infrared Spectroscopy

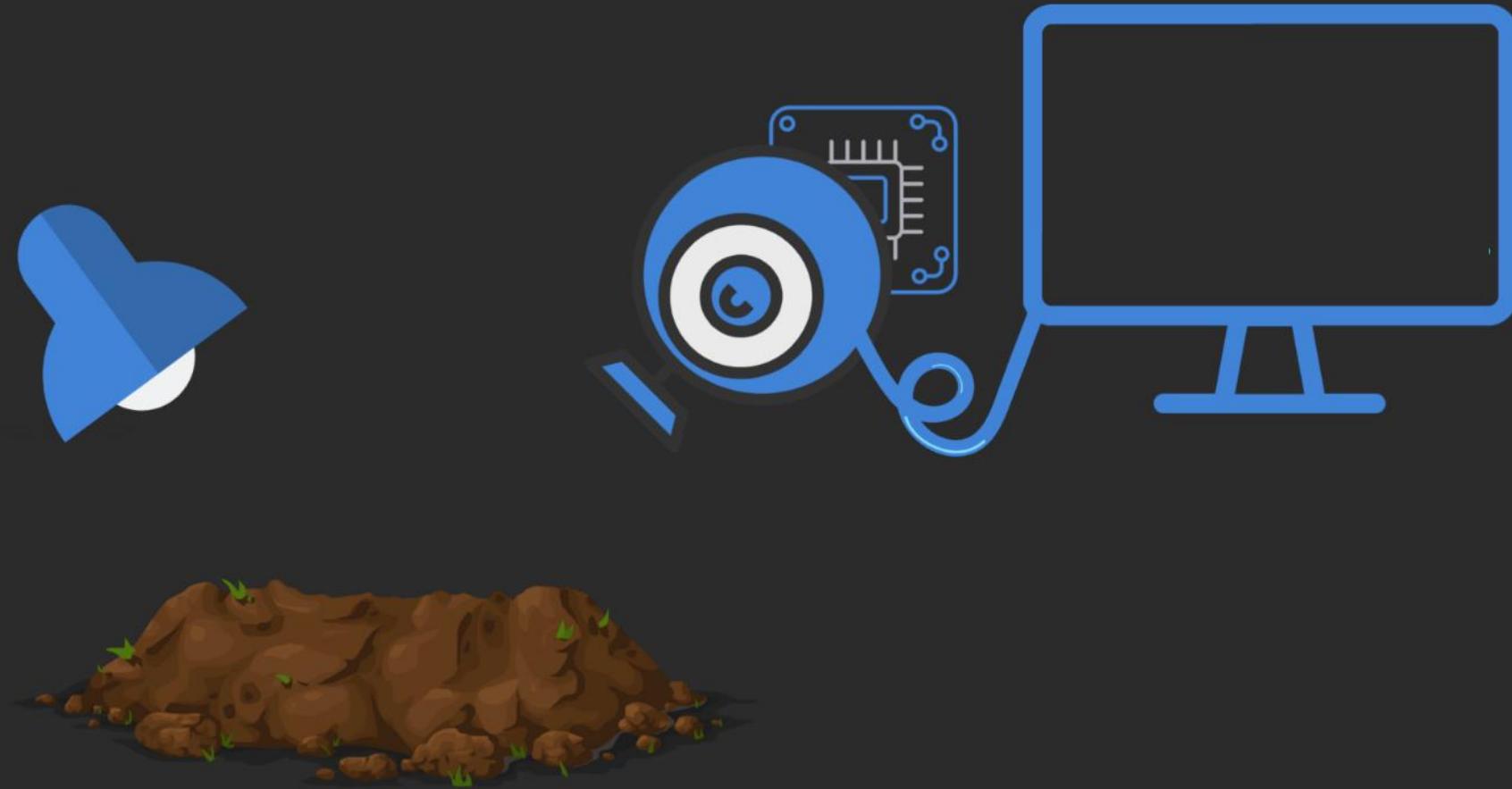
TropiRes – Summer School – 2024

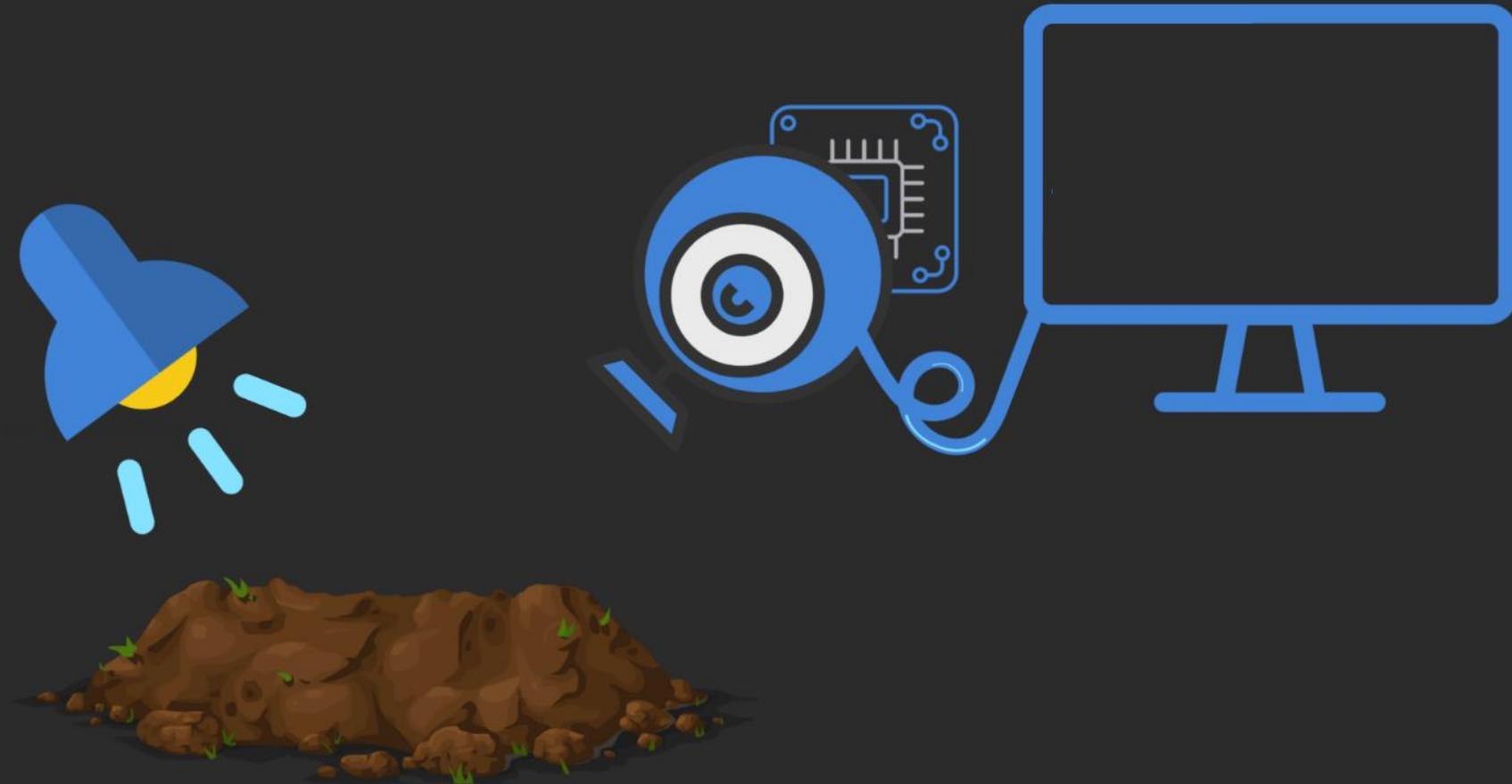
Module 2

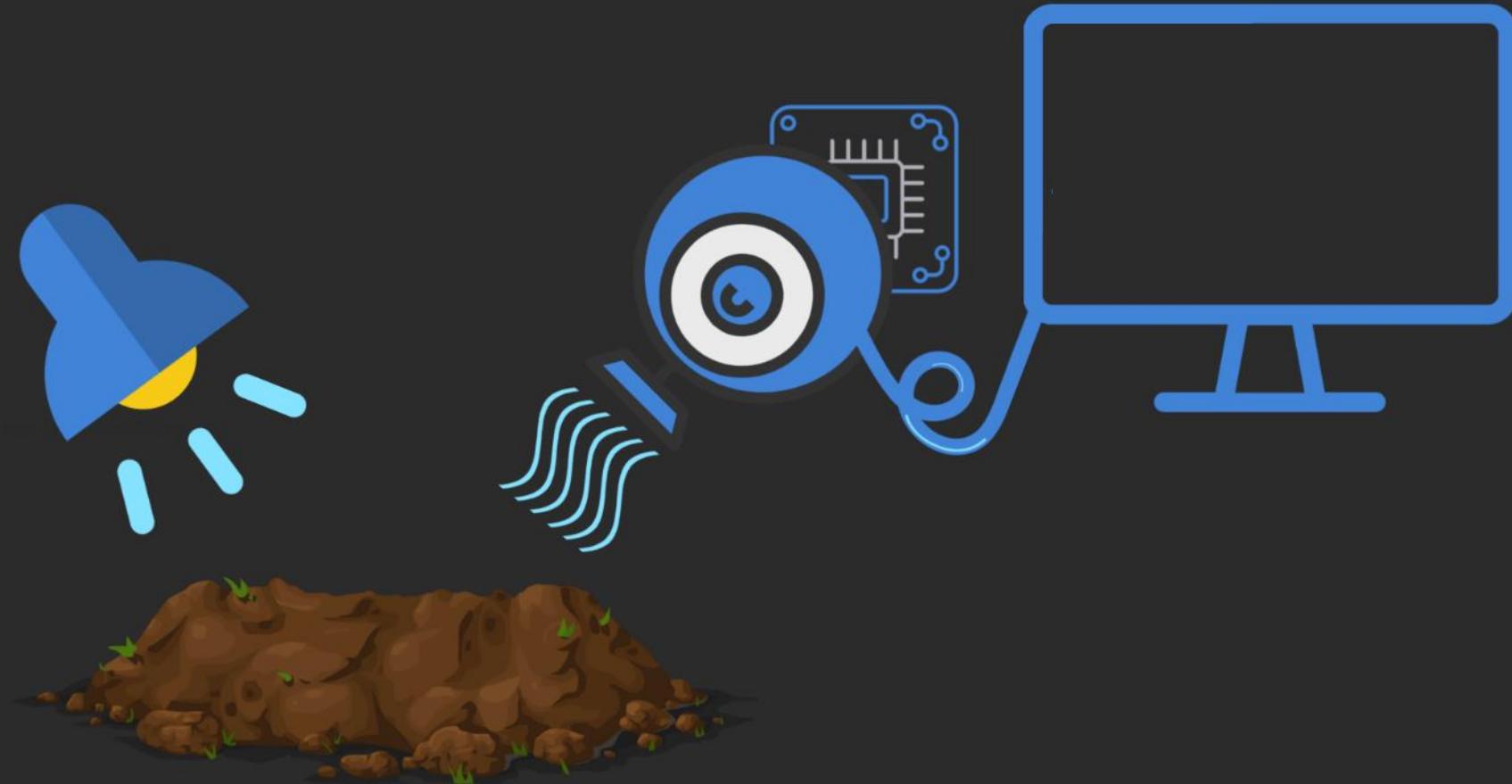
Leo Ramirez-López, Laura Summerauer, Moritz Mainka

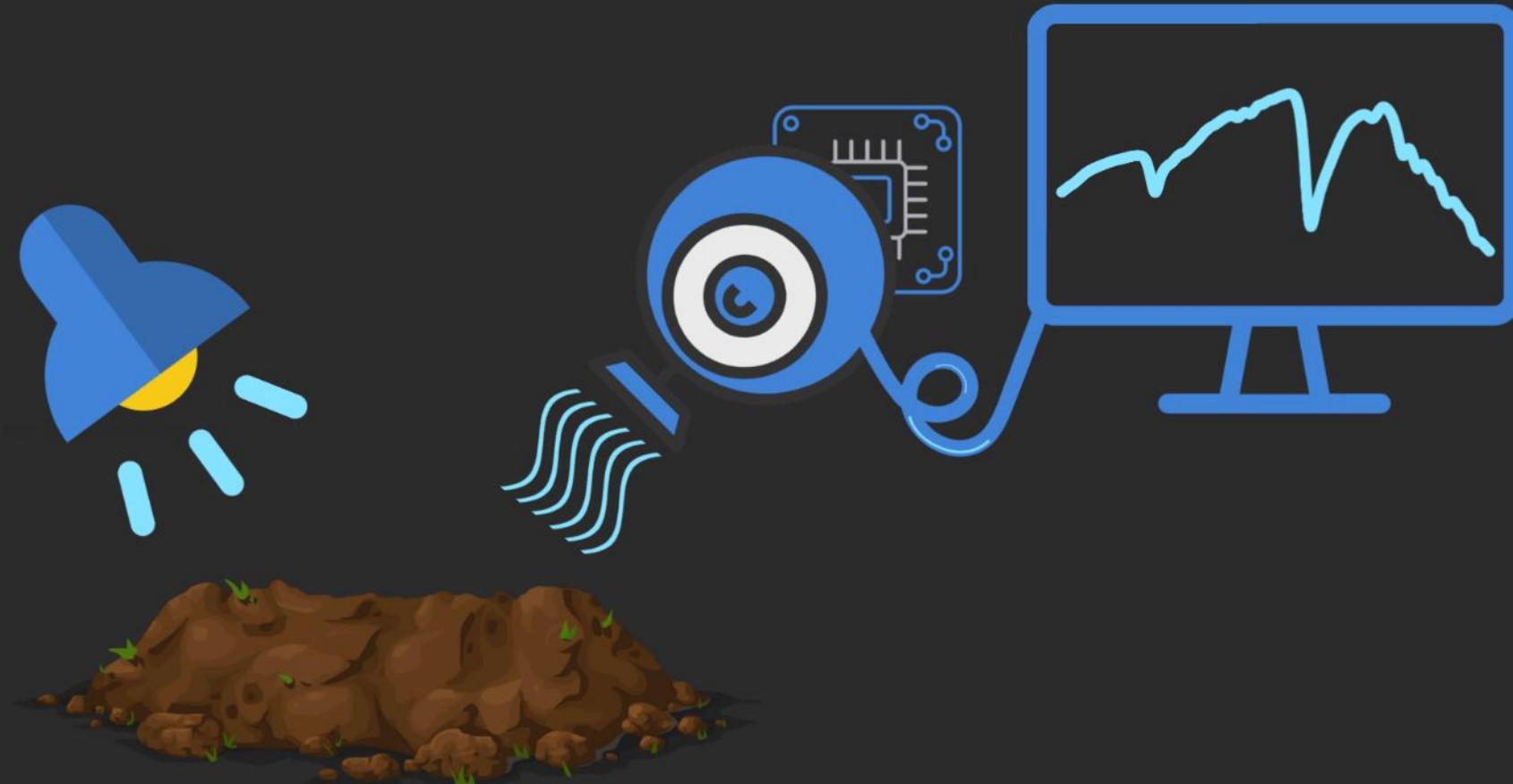
NIR/IR Spectroscopy – An equation for success

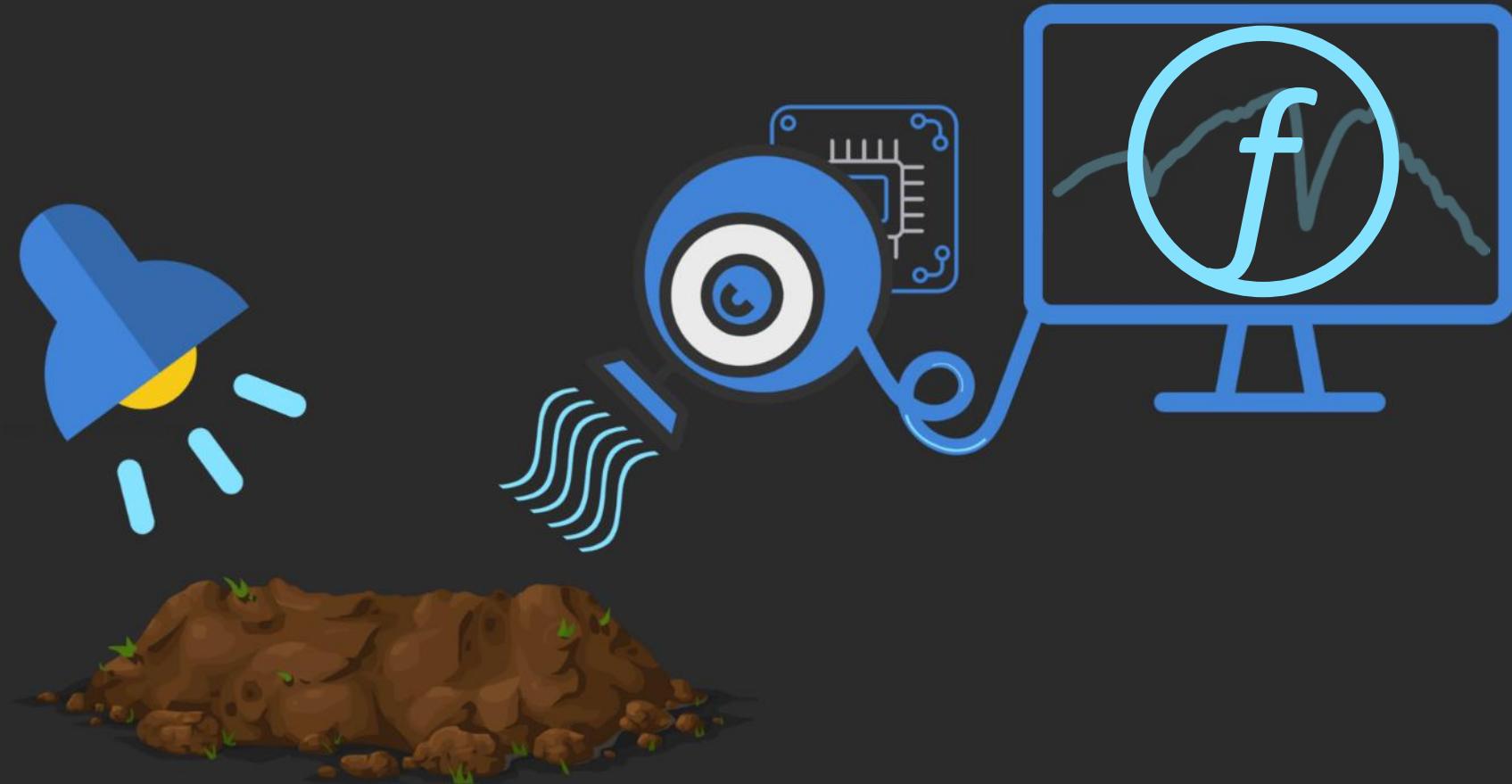
$$S_{\text{uccess}} = f(\text{hardware, software, data, models, people, automation, ...})$$

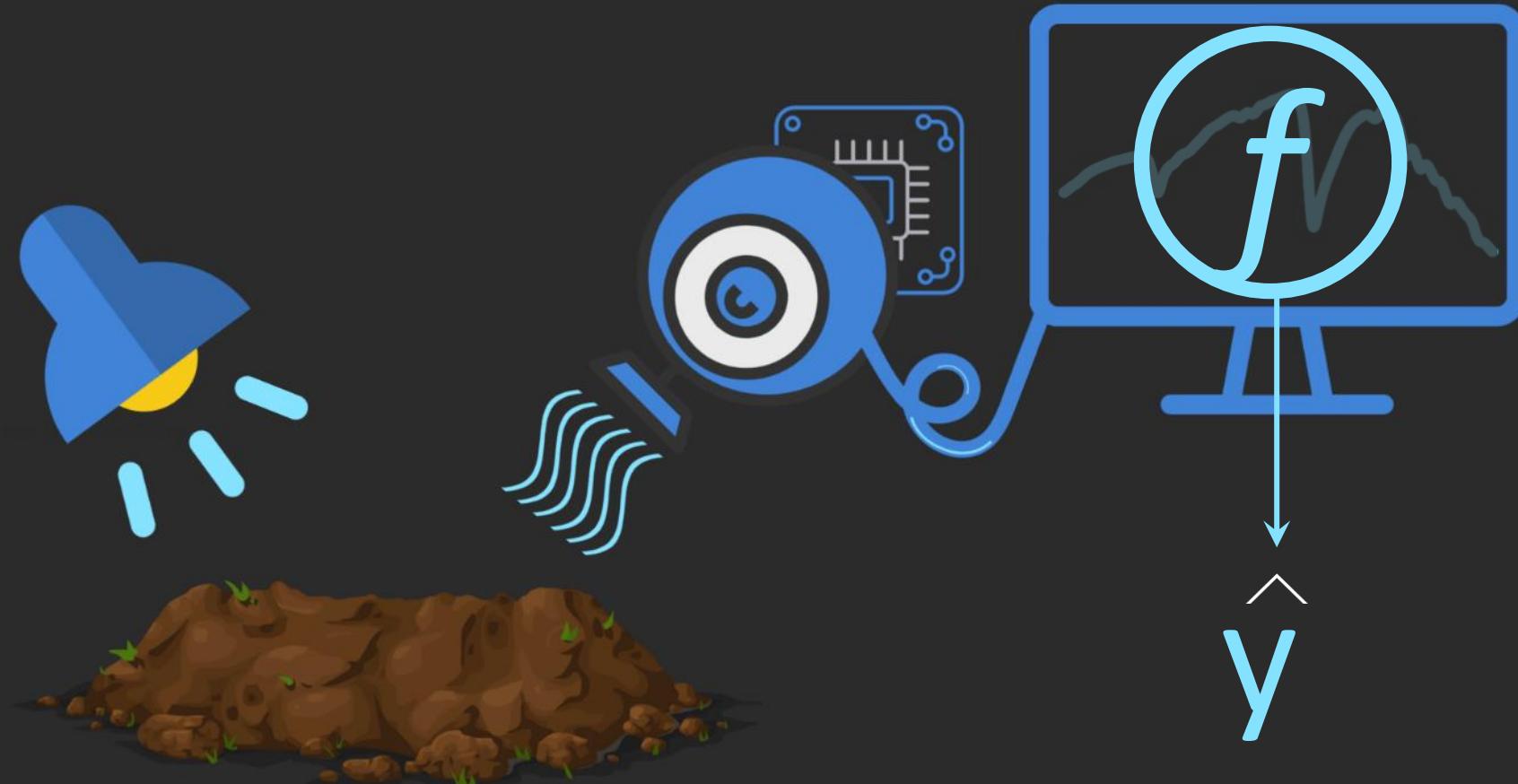






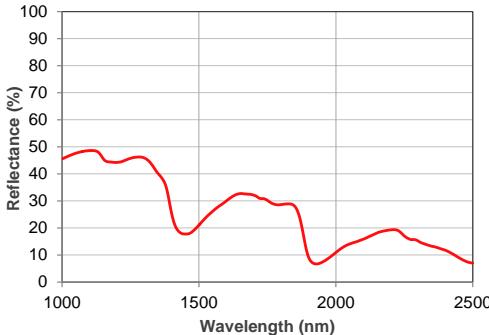






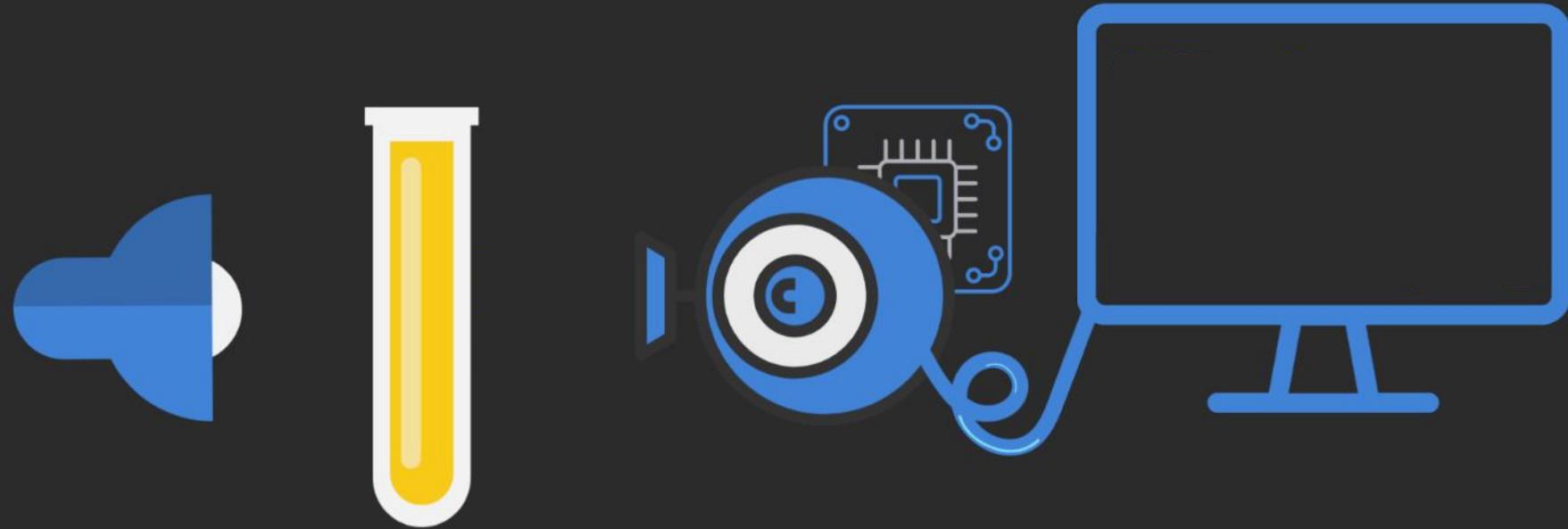
translating spectral information that we have into
chemical information that users need!

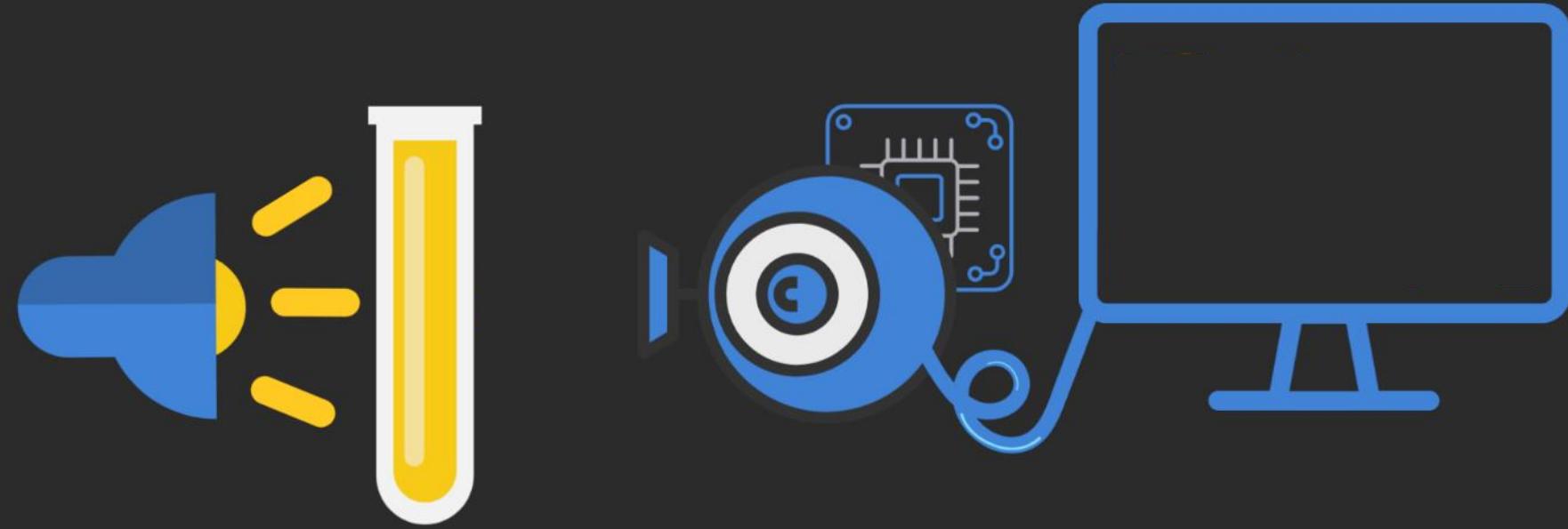
$$\hat{y} = f(\text{...}) + \epsilon$$

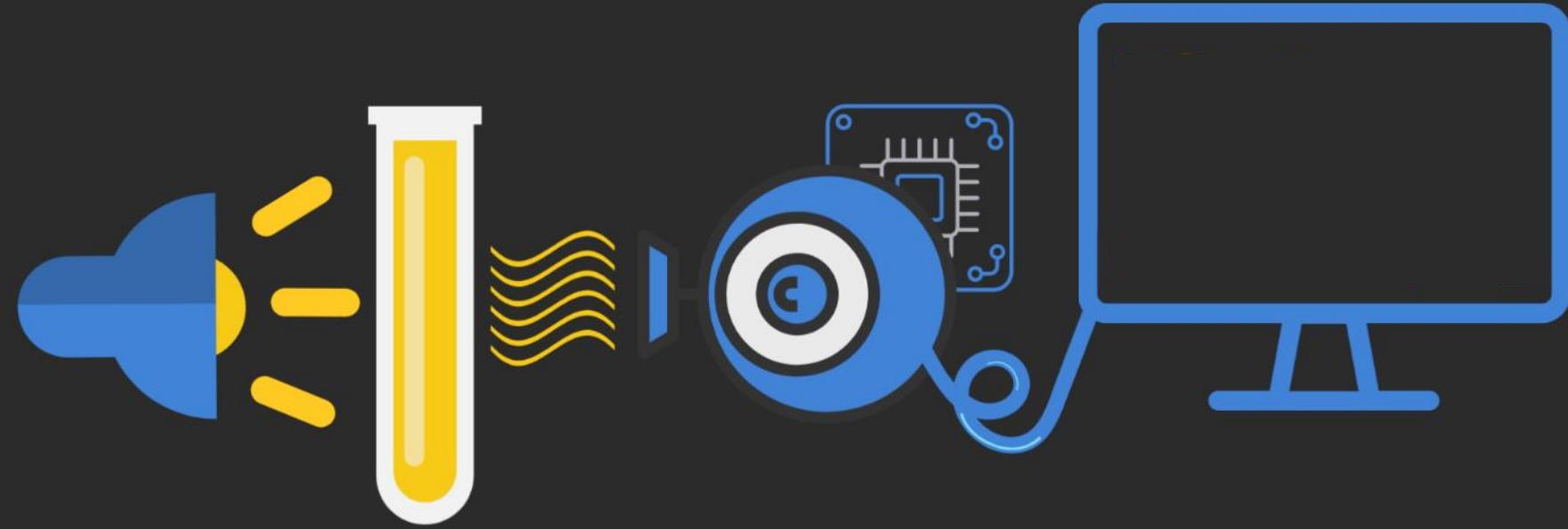


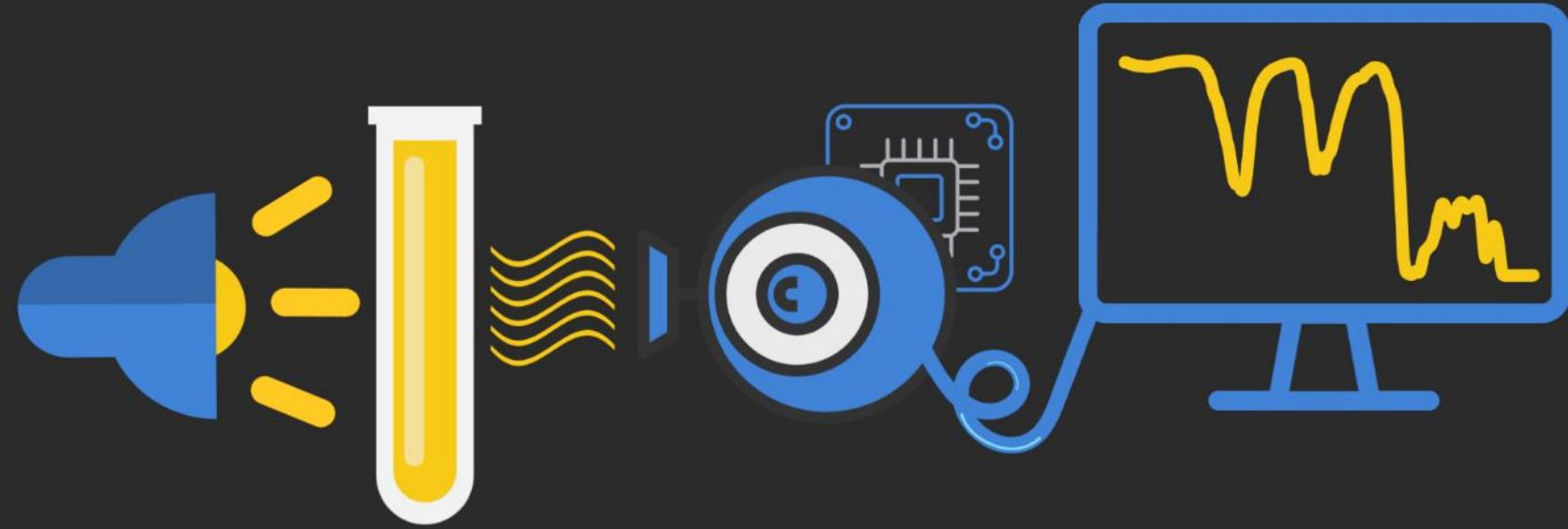
Function that can be
approximated by using
chemometrics/machine learning

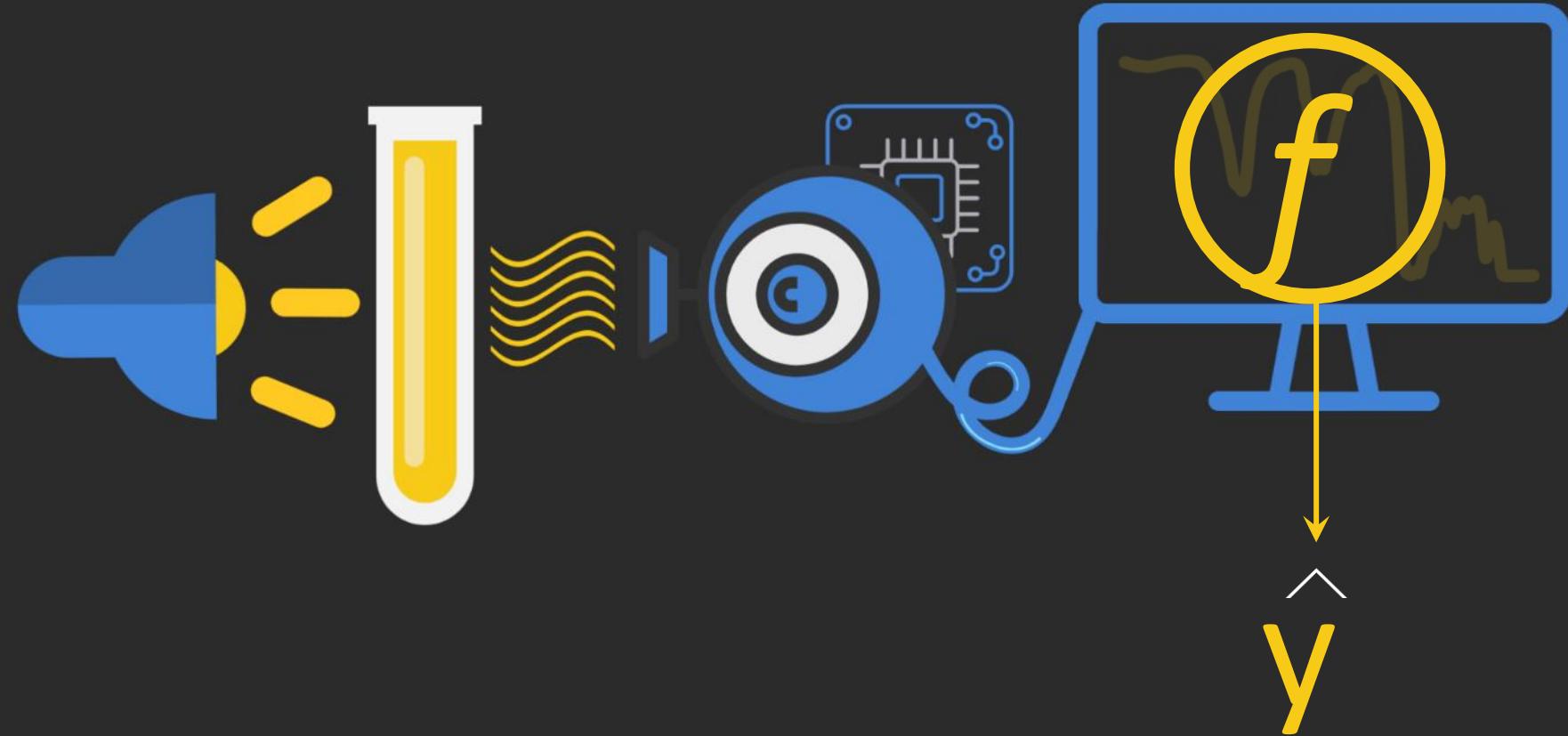






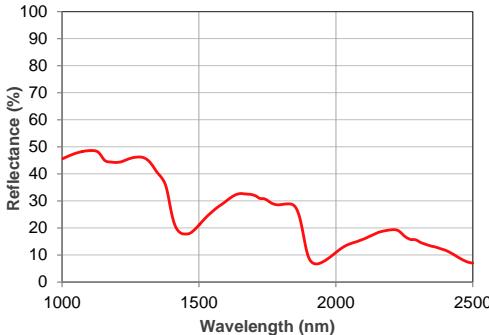






translating spectral information that we have into
chemical information that users need!

$$\hat{y} = f(\text{...}) + \epsilon$$



Function that can be
approximated by using
chemometrics/machine learning

transmission

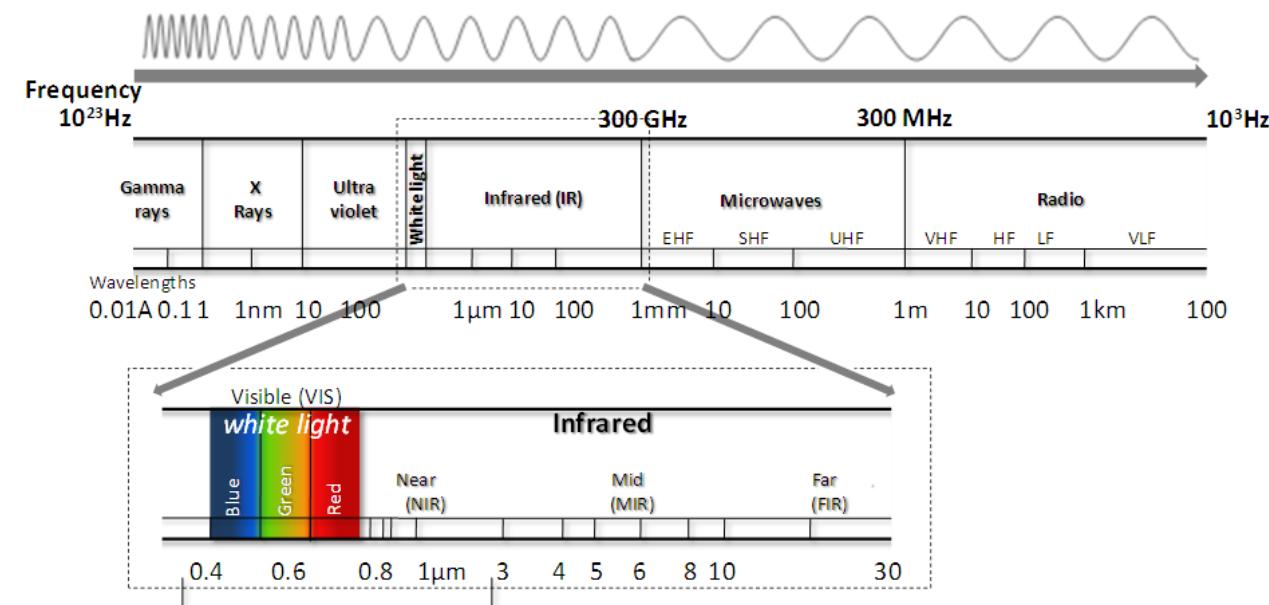
reflection

absorbance

transmission

reflection

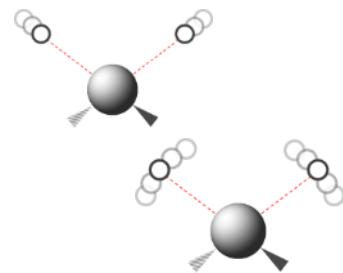
absorbance [$\log_{10}(1/\text{reflection})$]



NIR/IR spectroscopy is based on the measurements of the spectral response of objects in the infrared region of the electromagnetic spectrum.

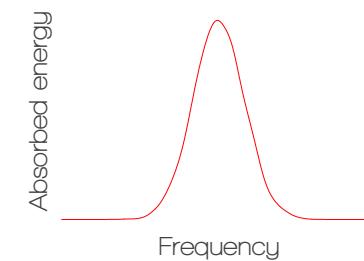
Near-infrared spectroscopy (NIR)...

How does it work?



A given molecule absorbs energy (light) if the frequency of the light energy matches the vibrational frequency of the functional group of the molecule.

The amount of energy absorbed can be different for each frequency (wavenumber). This usually results in energy absorption peaks.





Device name	Manufacturer	Spectral range (nm)	Price (USD)
<i>Bechtop</i>			
ASD FieldSpec*	Malvern Panalytical	350 – 2500	~68100
NIRS DS2500 L	FOSS Analytics	400 – 2500	~75000
<i>Handheld</i>			
SR-3500	Spectral Evolution	350 - 2500	~49700
MicroNIR 1700	VIAVI Solutions	950 - 1650	~20000
trinamiX	trinamiX/BASF	1450 - 2450	~8733
NeoSpectra	Si-Ware Systems	1300 - 2600	~3550
Agrocares	Agrocares	1300 - 2600	~6800
STS-VIS	Ocean Insight	350 - 830	~3408
NIRONE Sensor S	Spectral Engines	1750 - 2150	~3635
microPHAZIR	Thermo Fisher	1596 - 2396	?
SCiO	Consumer Physics	740 - 1070	~450
NIRscan	Texas Instruments	900 - 1700	~3000
Tellspec	Tellspec	900 - 1700	~320 - 1500

*a portable option is also available



ASD FieldSpec



FOSS



NeoSpectra



trinamix



Spectral Evolution



Agrocares

Device name	Manufacturer	Spectral range (nm)	Price (USD)
<i>Bechtop</i>			
ASD FieldSpec*	Malvern Panalytical	350 – 2500	~68100
NIRS DS2500 L	FOSS Analytics	400 – 2500	~75000
<i>Handheld</i>			
SR-3500	Spectral Evolution	350 - 2500	~49700
MicroNIR 1700	VIAVI Solutions	950 - 1650	~20000
trinamiX	trinamiX/BASF	1450 - 2450	~8733
NeoSpectra	Si-Ware Systems	1300 - 2600	~3550
Agrocares	Agrocares	1300 - 2600	~6800
STS-VIS	Ocean Insight	350 - 830	~3408
NIRONE Sensor S	Spectral Engines	1750 - 2150	~3635
microPHAZIR	Thermo Fisher	1596 - 2396	?
SCiO	Consumer Physics	740 - 1070	~450
NIRscan	Texas Instruments	900 - 1700	~3000
Tellspec	Tellspec	900 - 1700	~320 - 1500

*a portable option is also available



NeoSpectra



trinamix

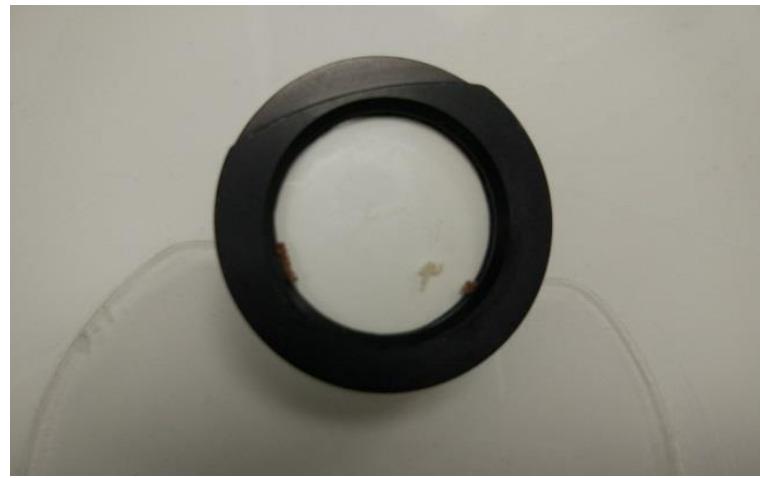


Spectral Evolution



Photometric alignment/calibration





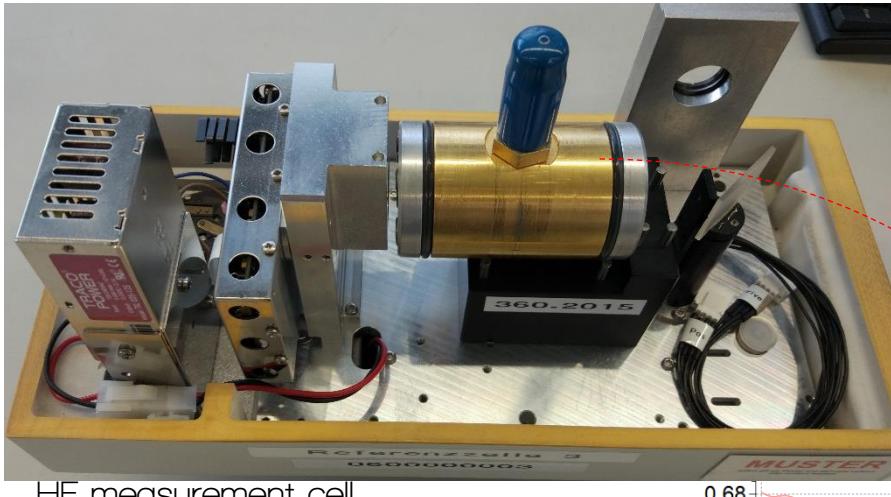
Photometric alignment/calibration



Wavelength alignment

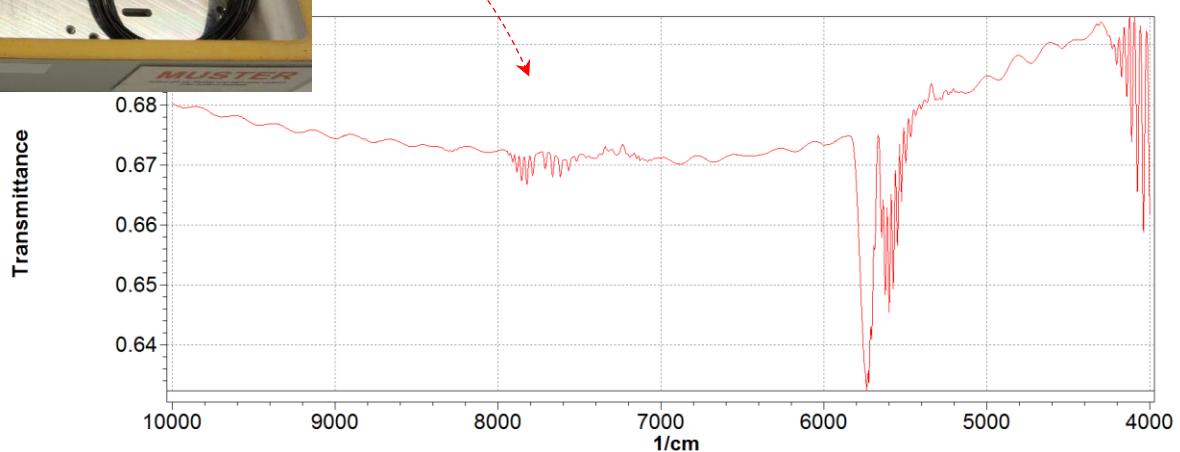
1. Hydrogen fluoride alignment

Some instruments are factory standardized using a wavelength standardization cell containing Hydrogen Fluoride (HF) gas. The spectral signature of this gas consist in very well defined peaks which can be measured with high precision and accuracy.



HF measurement cell

HF spectrum for wavelength alignment



Let's scan our first sample

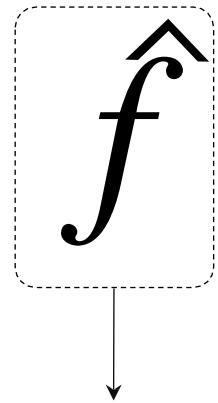


Download Neo-Spectra COLLECT
[Download Link GooglePlay](#)

Scanning in the IR spectrometer...



NIR Applications

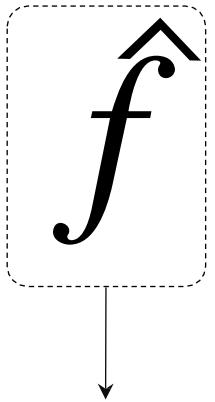


Function that can be
approximated by using
chemometrics/machine learning

We currently use our own BUCHI
software for that!

These functions are device-, product- and property-specific.

NIR Applications

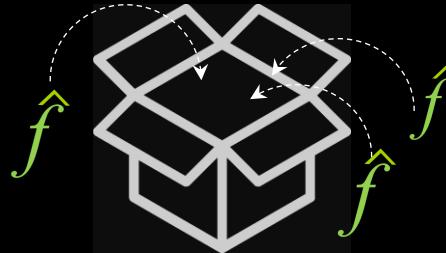


Function that can be approximated by using chemometrics/machine learning

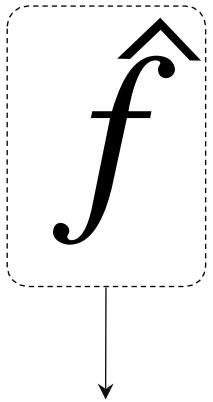
We currently use our own BUCHI software for that!

These functions are device-, product- and property-specific.

For each product we build and group a set of functions that quantify relevant product properties (e.g. moisture, protein, fat). This group of functions is called application.



NIR Applications

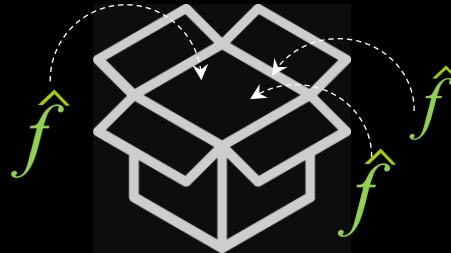


Function that can be approximated by using chemometrics/machine learning

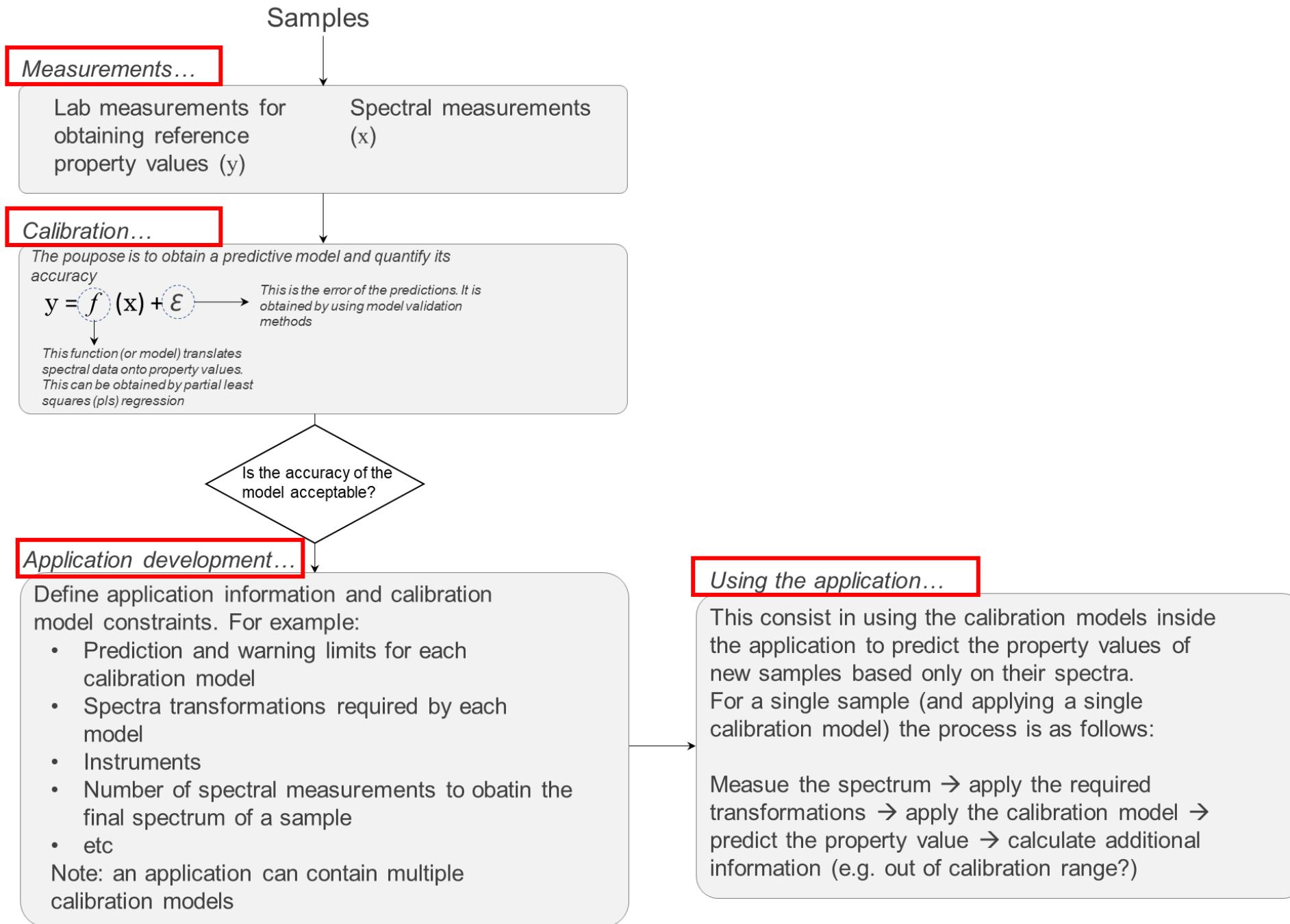
We currently use our own BUCHI software for that!

These functions are device-, product- and property-specific.

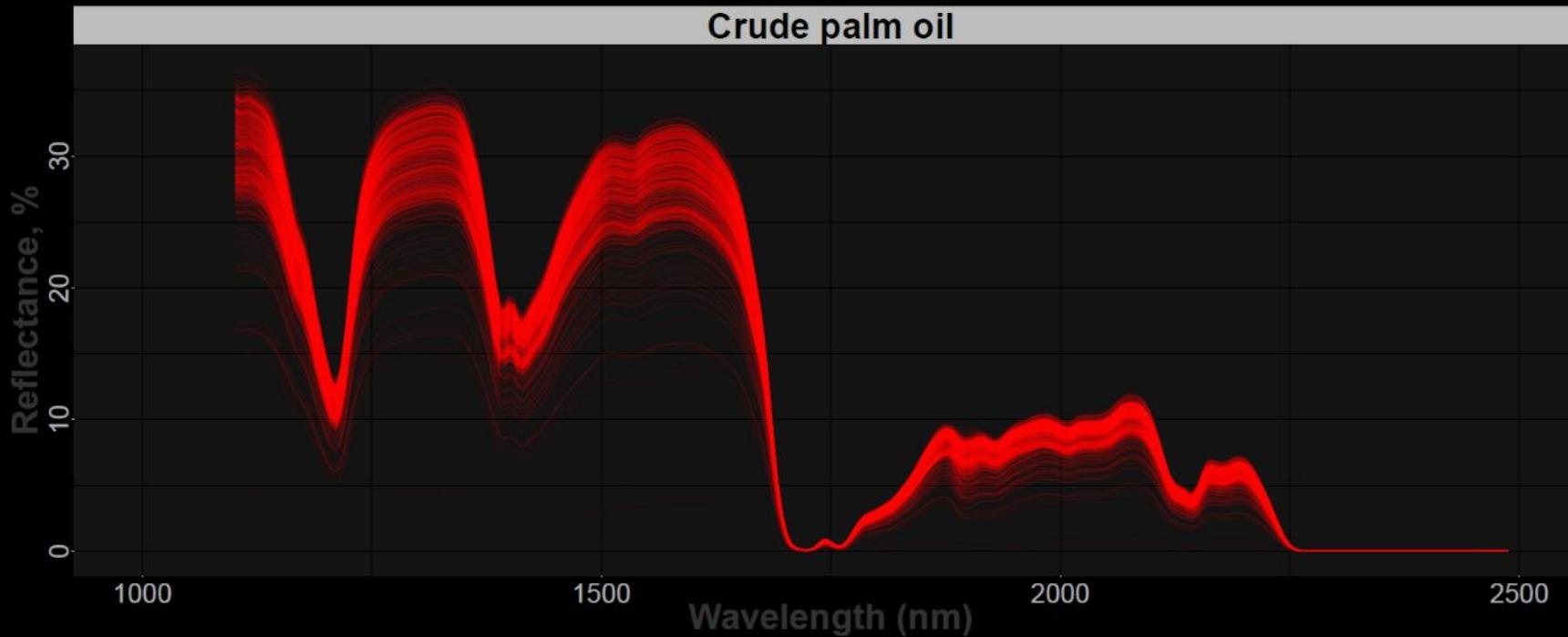
For each product we build and group a set of functions that quantify relevant product properties (e.g. moisture, protein, fat). This group of functions is called application.



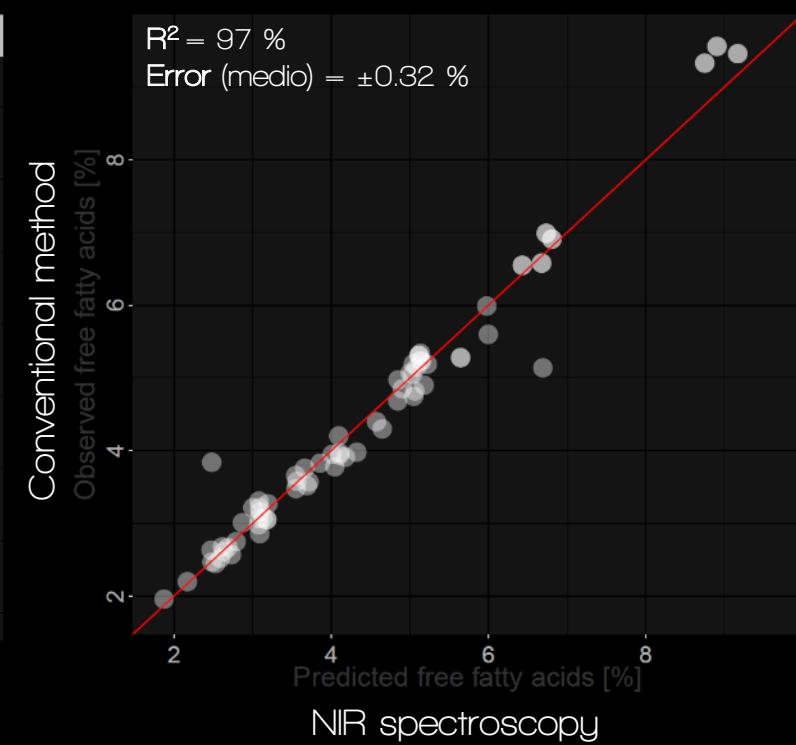
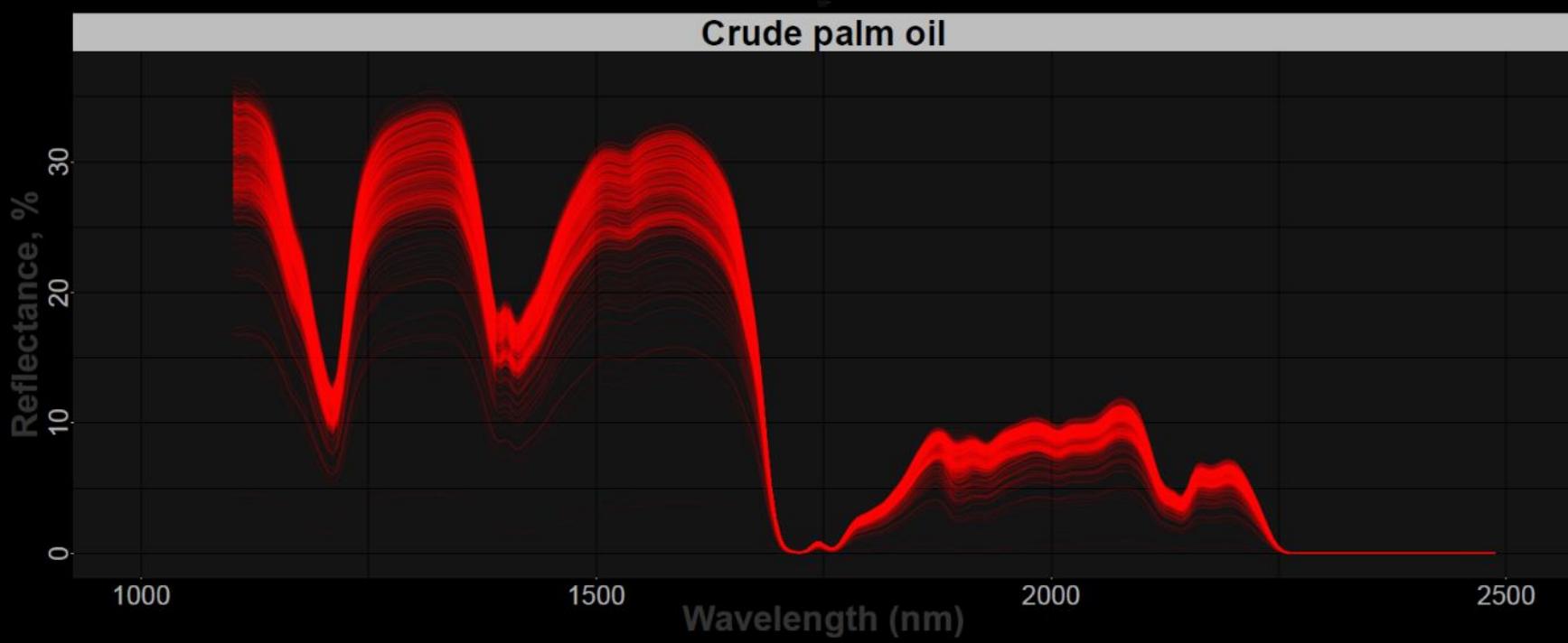
From average user's perspective, a NIR device without an application can be perceived as useless.



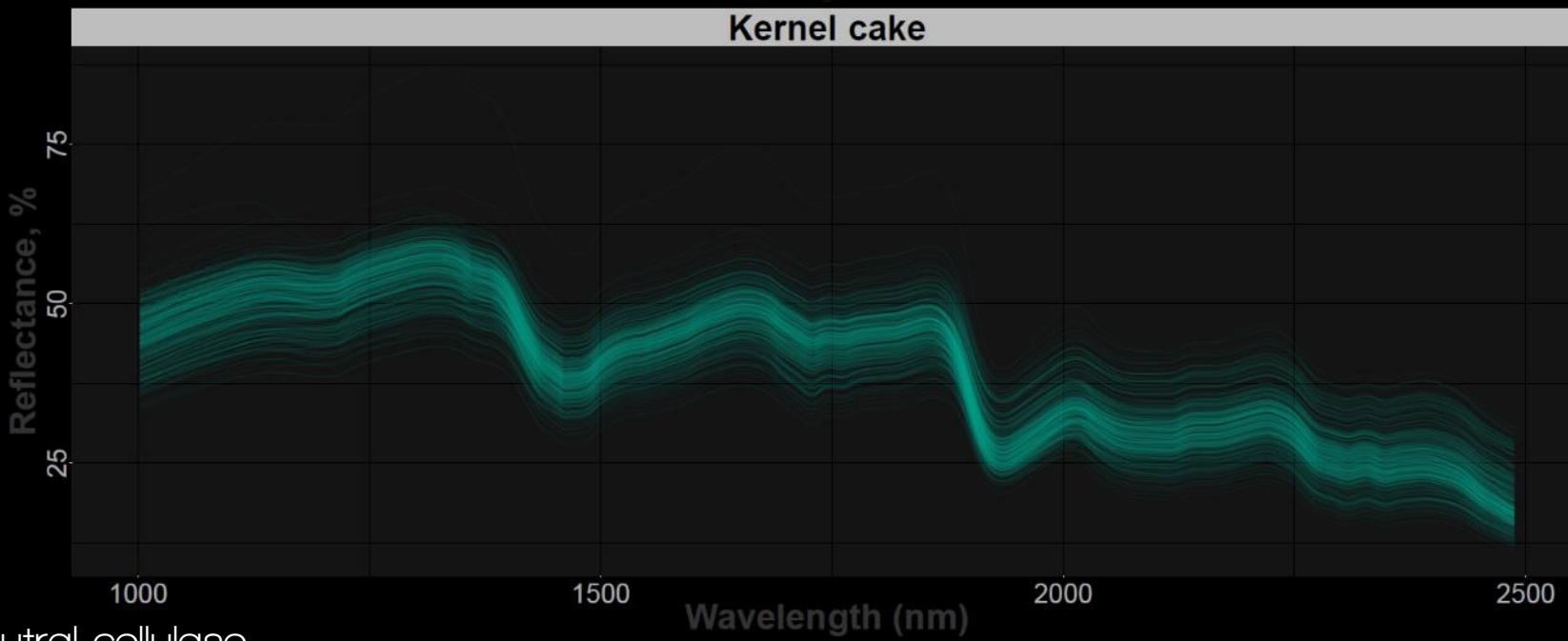
Crude palm oil



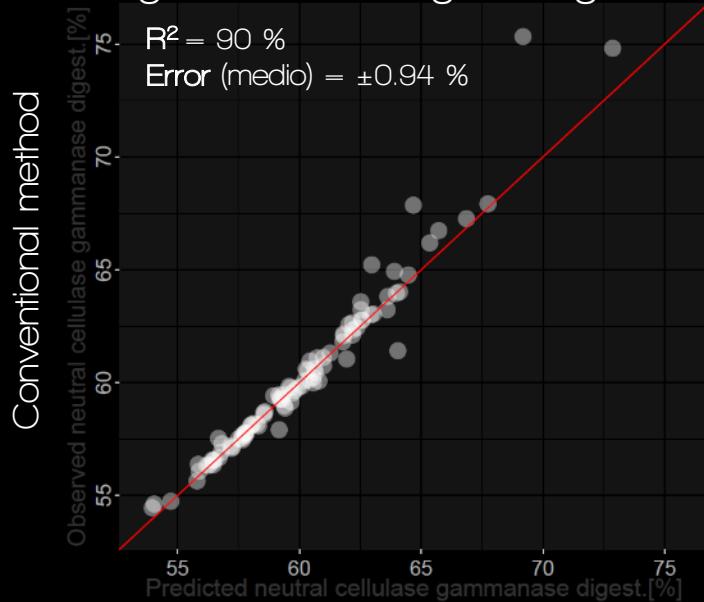
Free fatty acids



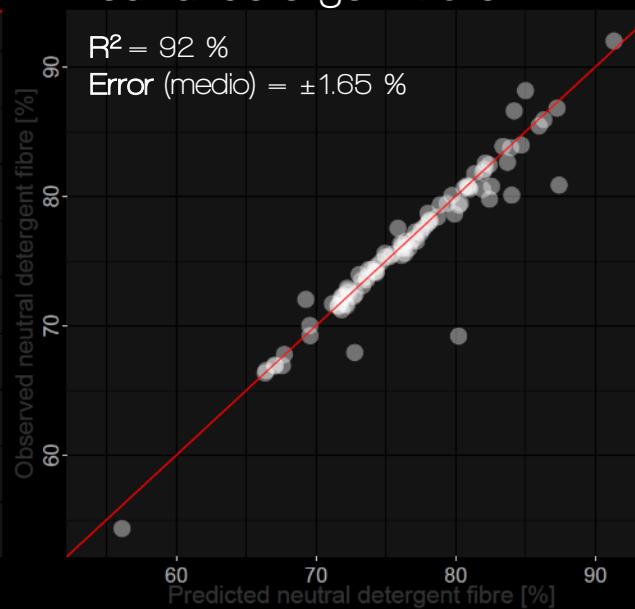
Kernel cake



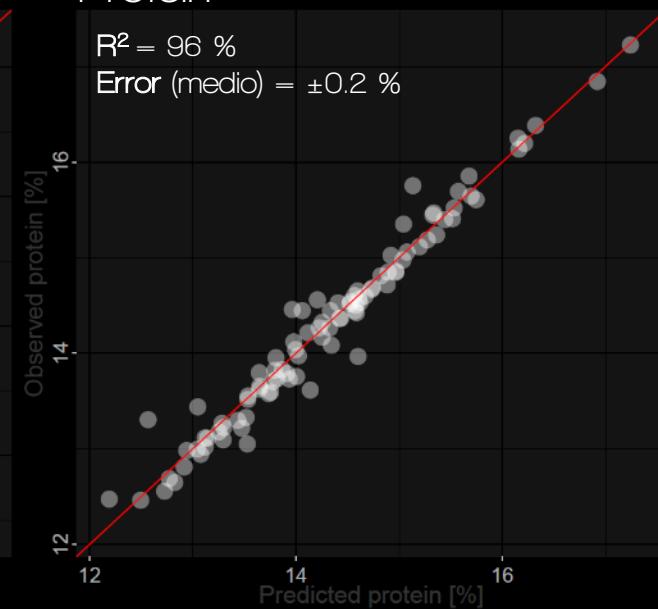
Neutral cellulase
gammanase digestibility



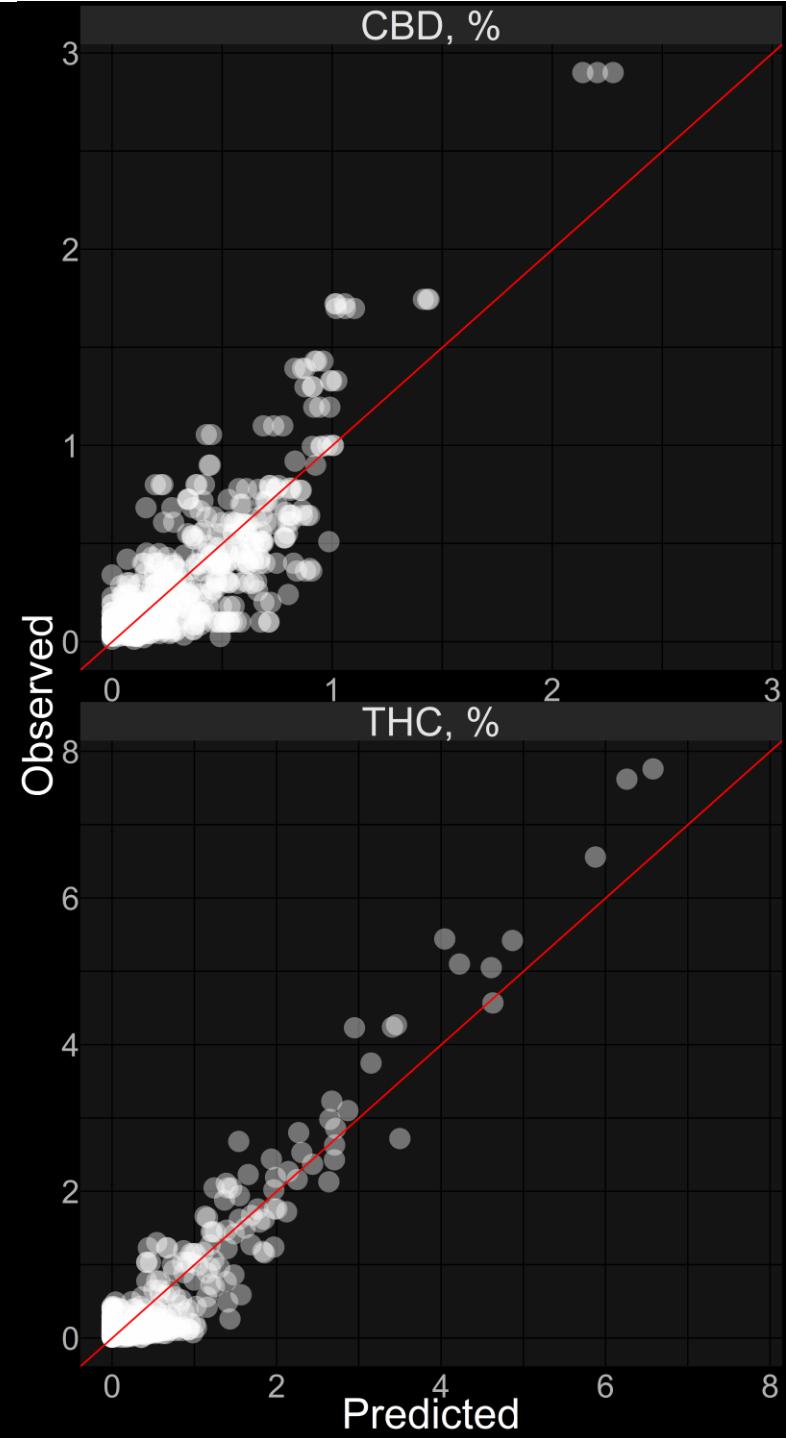
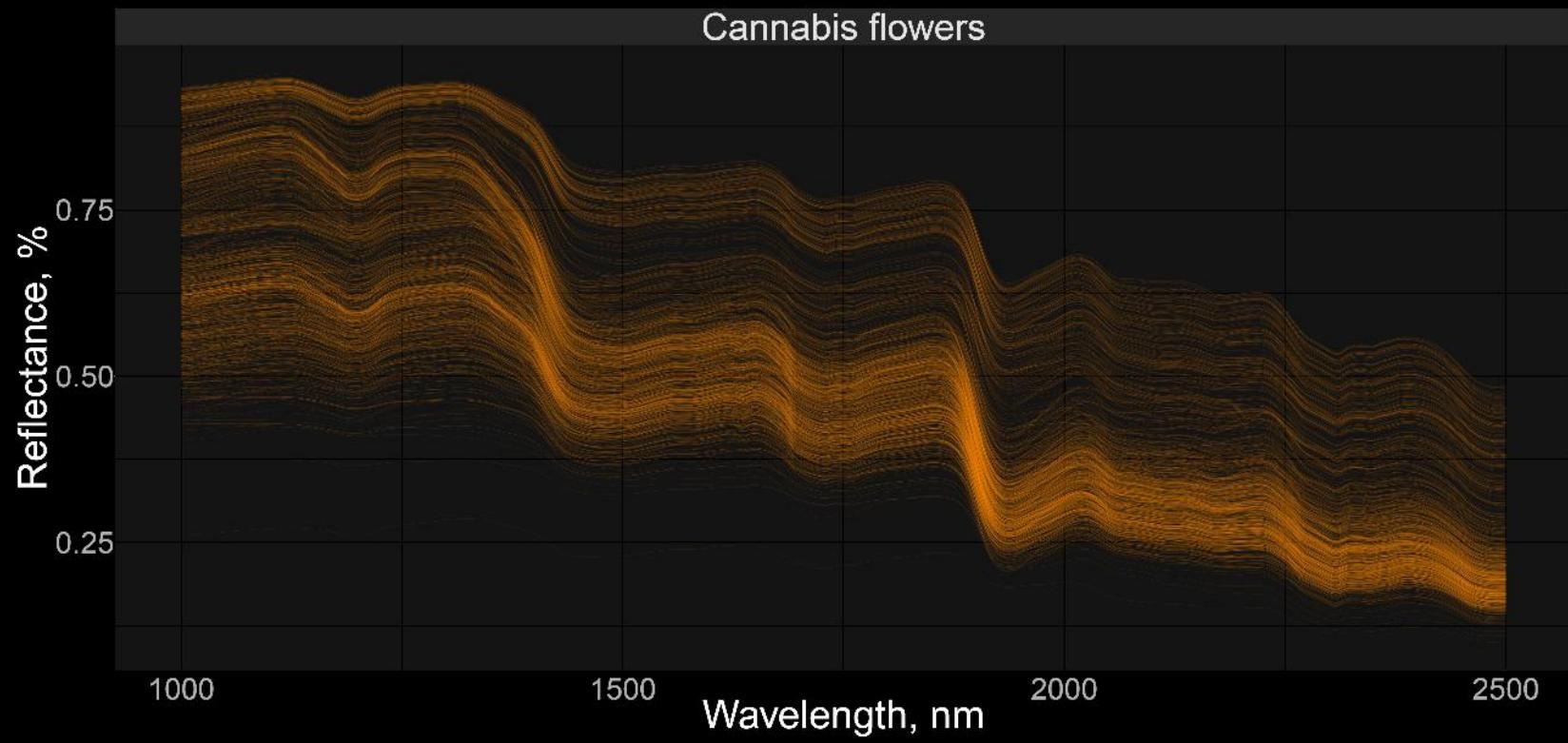
Neutral detergent fibre

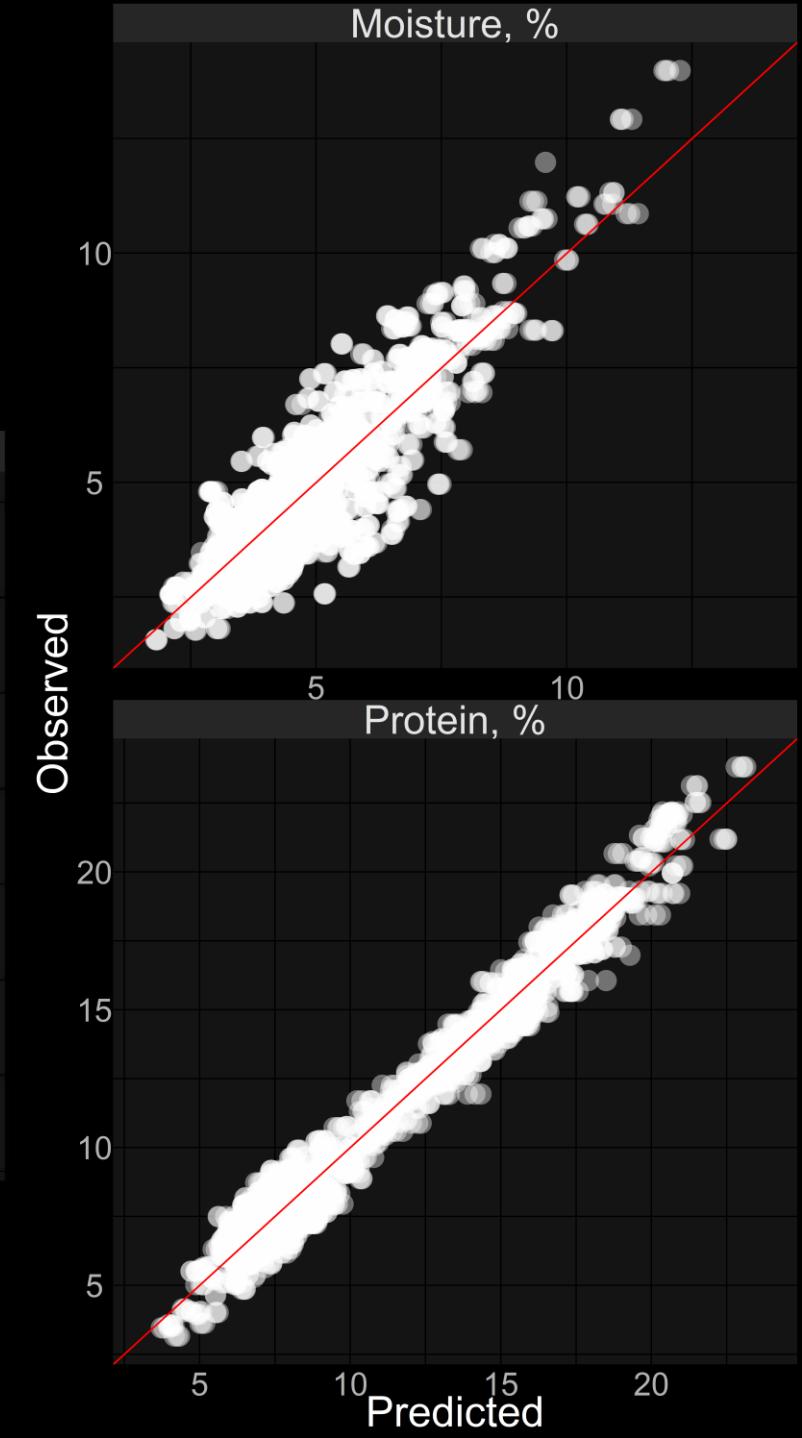
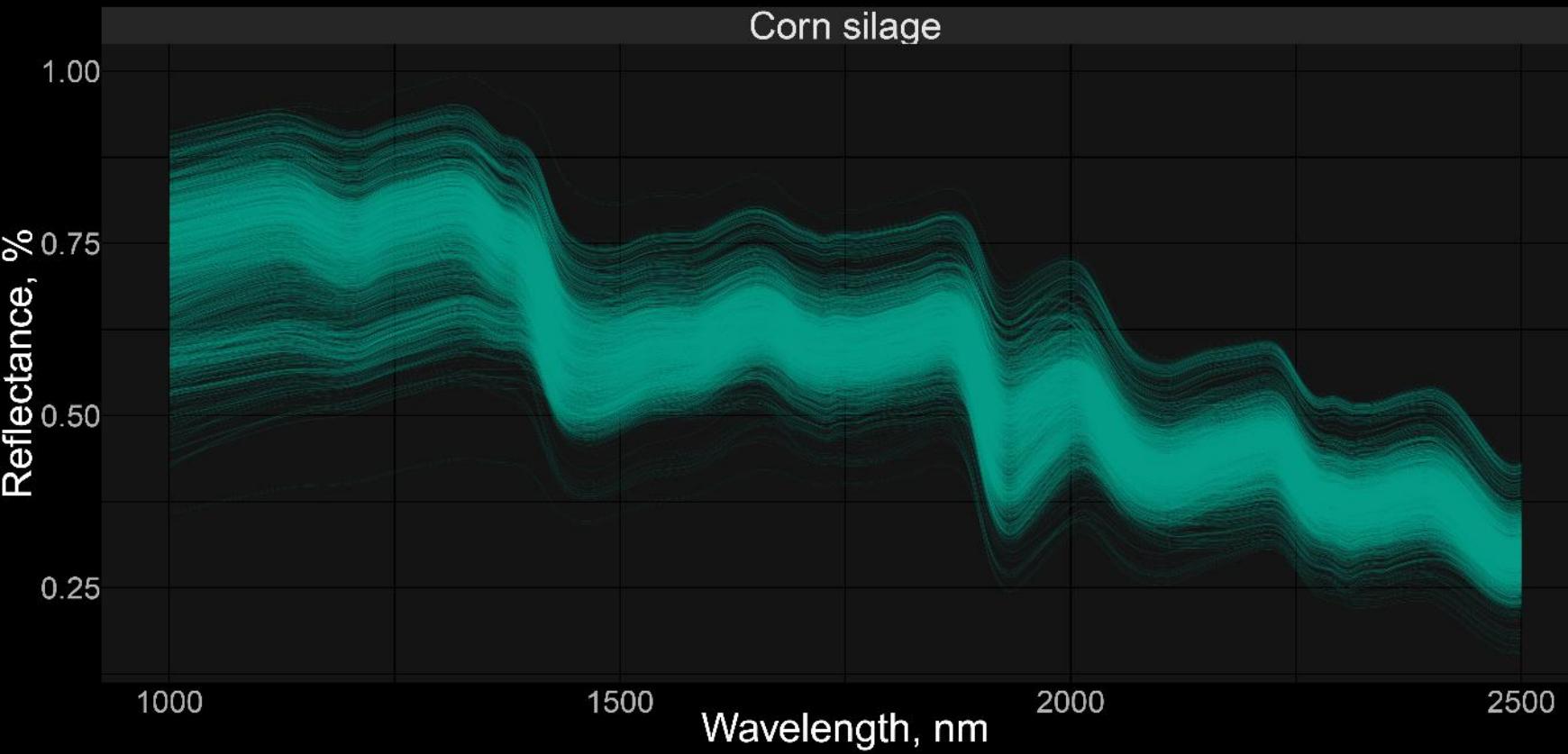


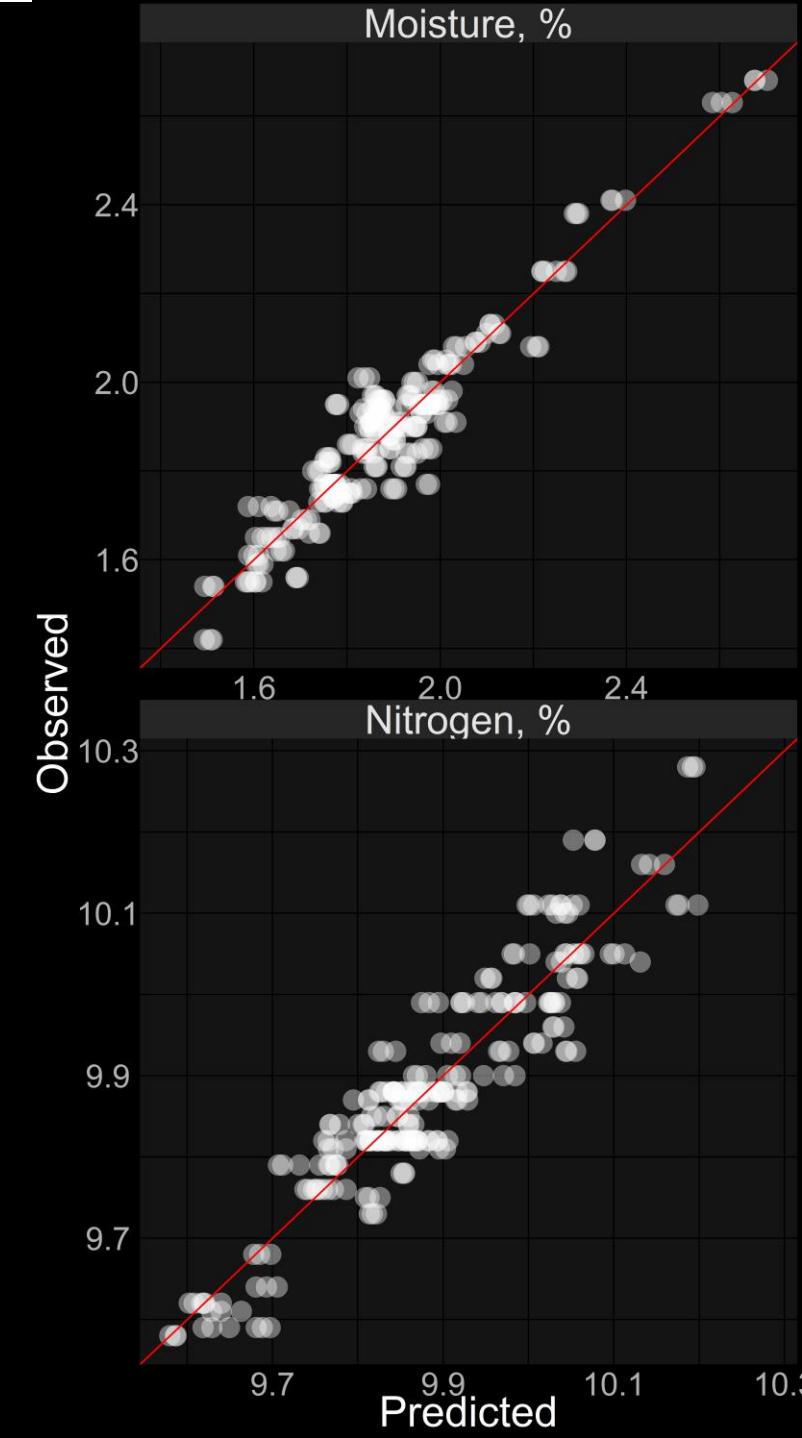
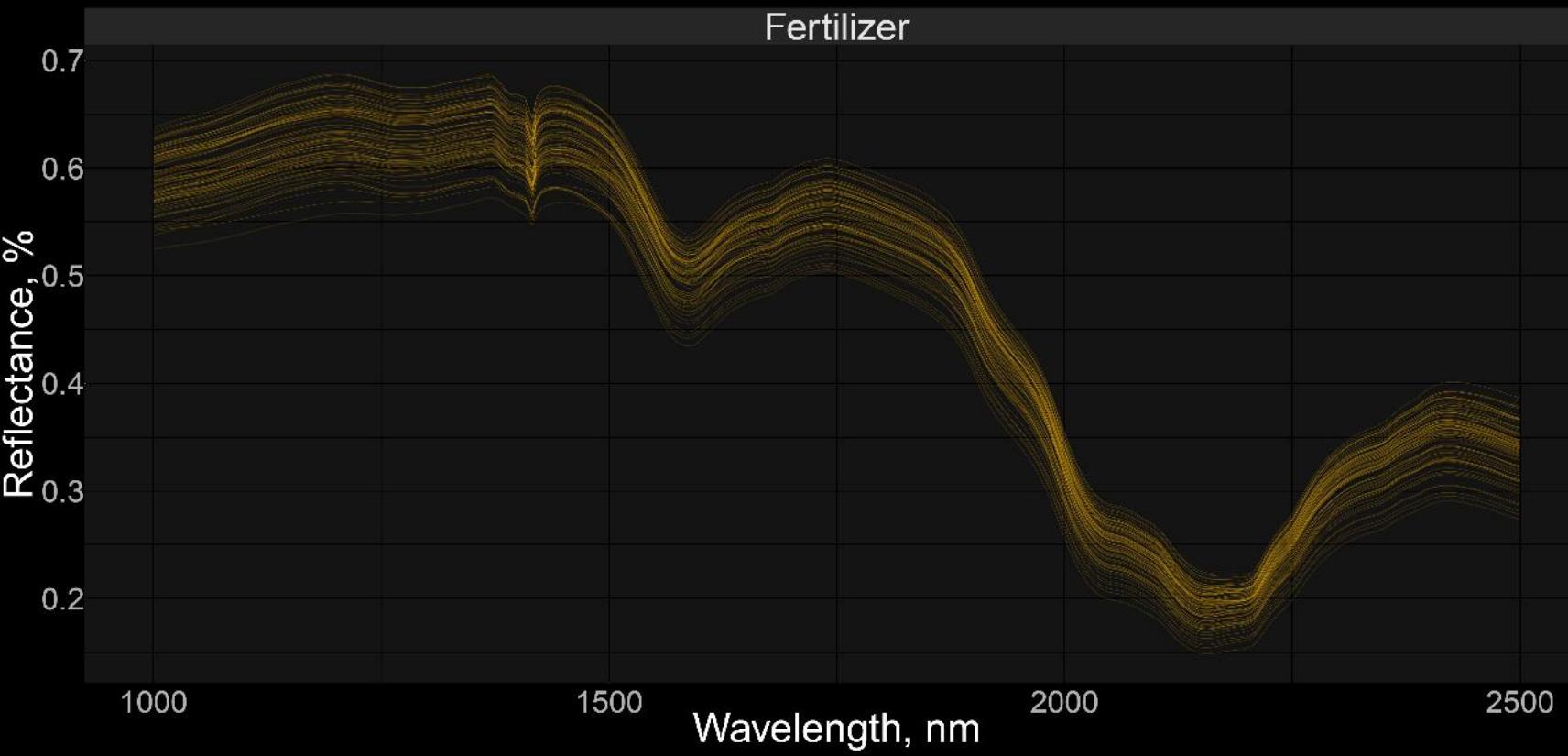
Protein

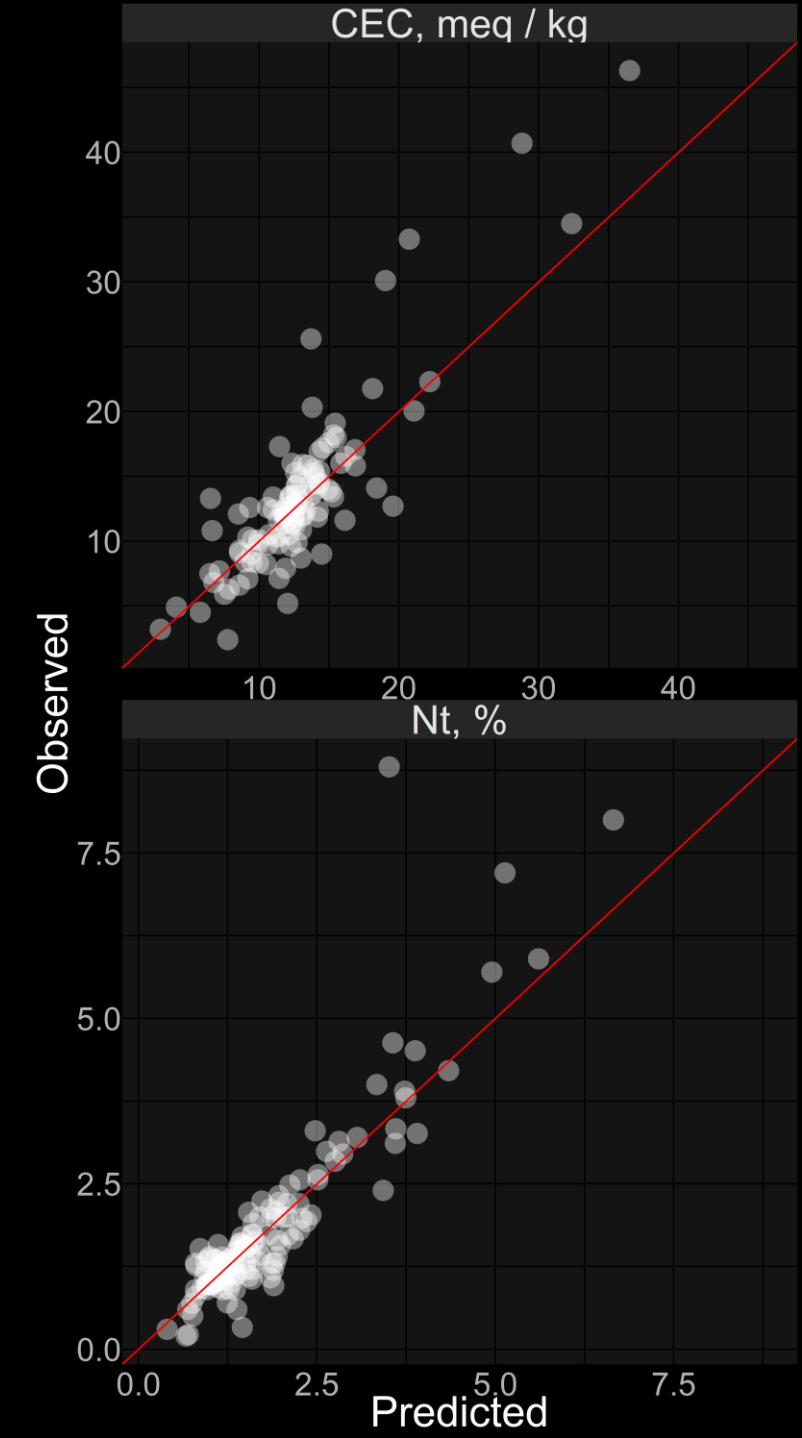
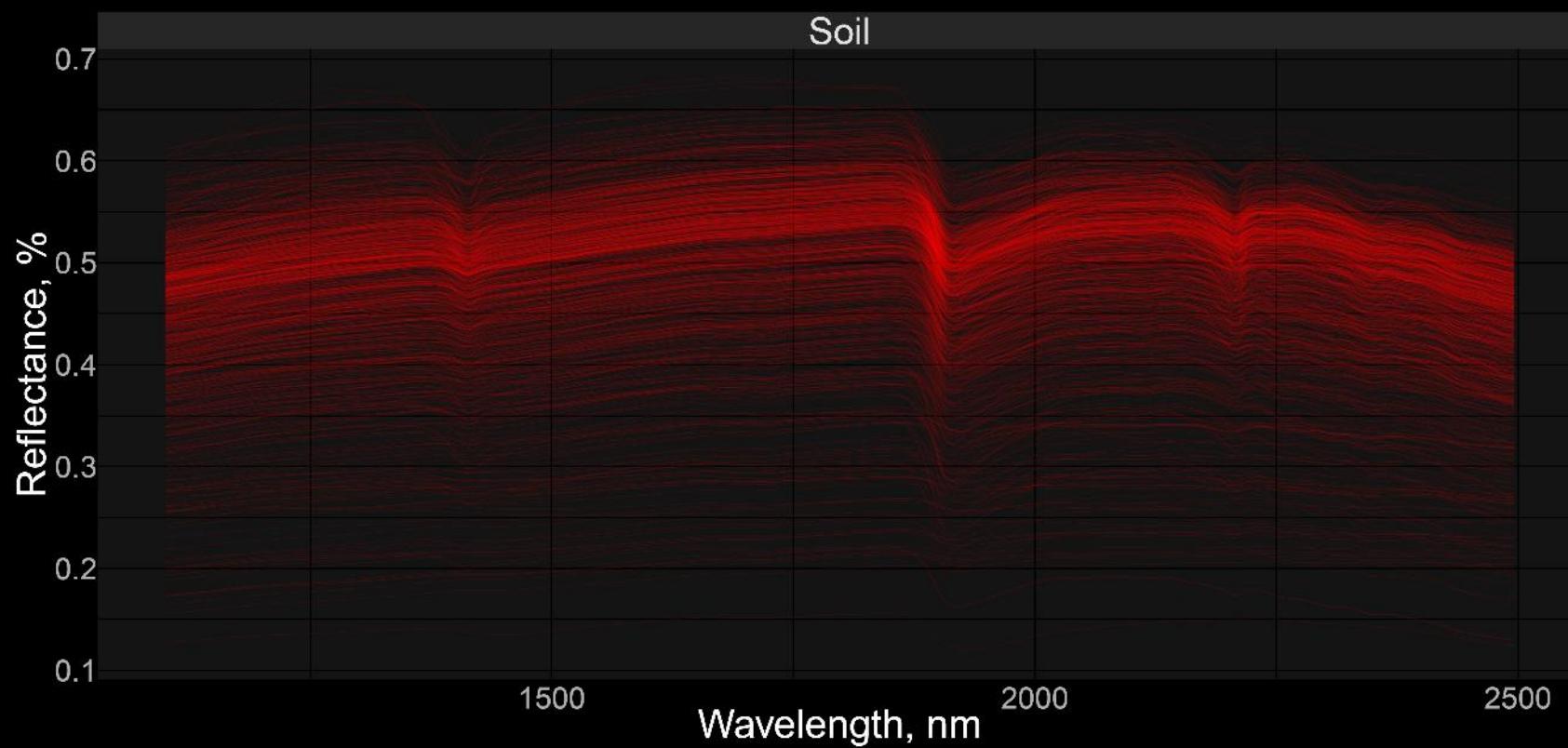
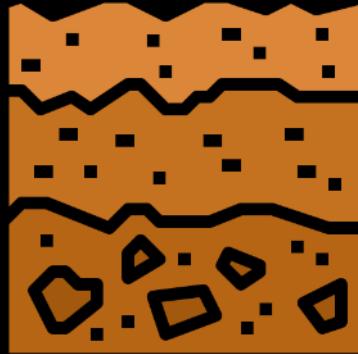


NIR spectroscopy









NIR/IR Spectroscopy – An equation for success

$$S_{\text{uccess}} = f(\text{hardware, software, data, models, people, automation, ...})$$

Application Development Roadmap



Collect calibration data

Legacy and/or new data to develop calibrations is collected.

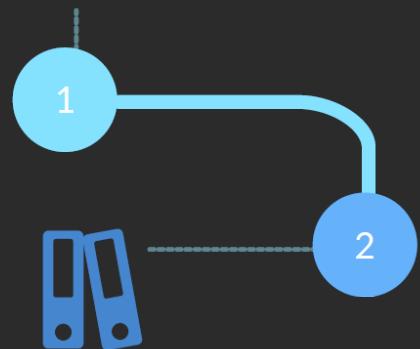
1

Application Development Roadmap



Collect calibration data

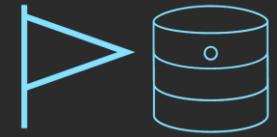
Legacy and/or new data to develop calibrations is collected.



Organize data

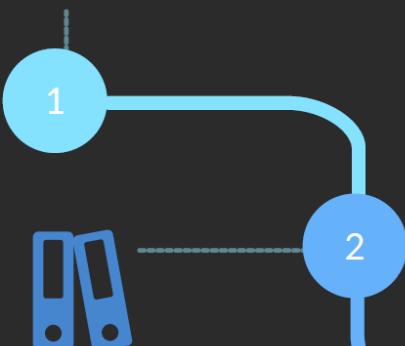
cleaning/filtering, sorting, harmonizing, pooling, storing and document all the calibration data gathered.

Application Development Roadmap



Collect calibration data

Legacy and/or new data to develop calibrations is collected.



Organize data

cleaning/filtering, sorting, harmonizing, pooling, storing and document all the calibration data gathered.



Calibrate spectral models

Use the calibration data to calibrate the models of the multiple properties of the target product. For each model, several data processing methods are empirically tested to optimize the accuracy and precision of the models.

The training/calibration set

Property
(response)

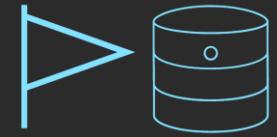
NIR signals (predictors)

n

n

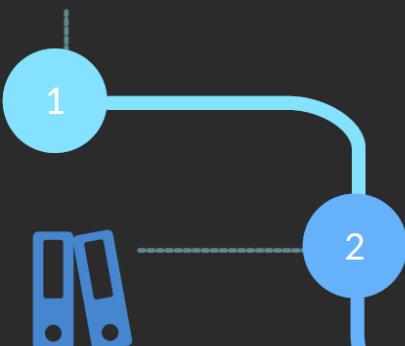


Application Development Roadmap



Collect calibration data

Legacy and/or new data to develop calibrations is collected.



Organize data

cleaning/filtering, sorting, harmonizing, pooling, storing and document all the calibration data gathered.

Calibrate spectral models

Use the calibration data to calibrate the models of the multiple properties of the target product. For each model, several data processing methods are empirically tested to optimize the accuracy and precision of the models.



The training/calibration set

Property
(response)



n

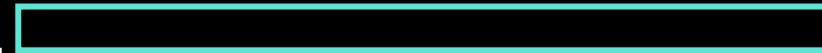
NIR signals (predictors)



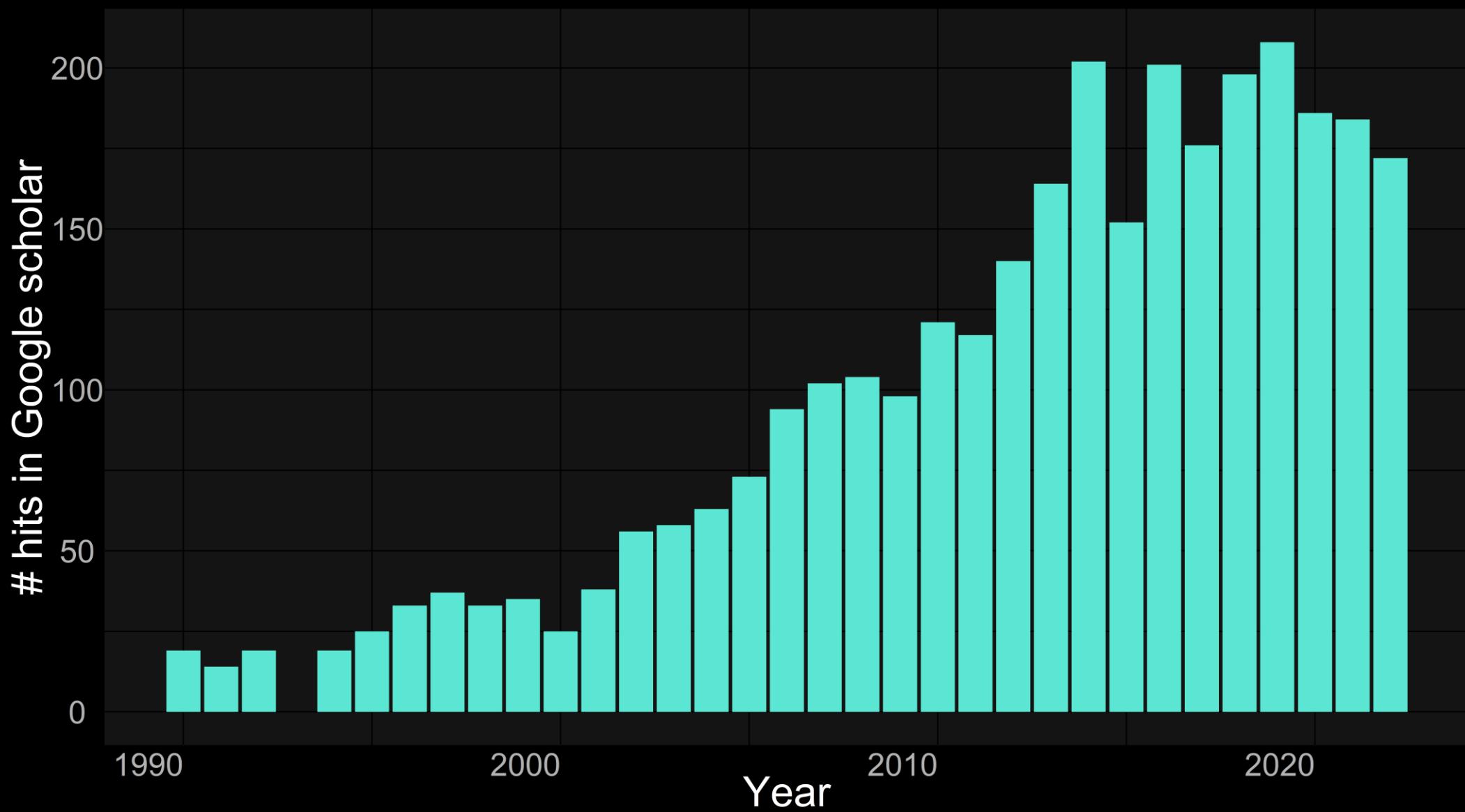
n

The calibration model \hat{f}

1

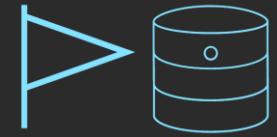


33 years
~3166 scientific publications



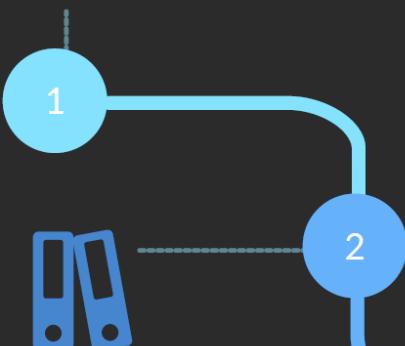
Search terms:
Soil + spectroscopy (in the title of the paper)

Application Development Roadmap



Collect calibration data

Legacy and/or new data to develop calibrations is collected.



Organize data

cleaning/filtering, sorting, harmonizing, pooling, storing and document all the calibration data gathered.



Calibrate spectral models

Use the calibration data to calibrate the models of the multiple properties of the target product. For each model, several data processing methods are empirically tested to optimize the accuracy and precision of the models.



and then?

Application Development Roadmap



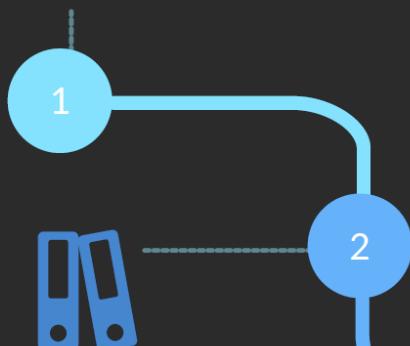
Pack into a predictive application

All the spectral models are packed into an application. All the technical and user-required documentation is generated along with the necessary licensing information.



Collect calibration data

Legacy and/or new data to develop calibrations is collected.



Organize data

Cleaning/filtering, sorting, harmonizing, pooling, storing and document all the calibration data gathered.

1



Calibrate spectral models

Use the calibration data to calibrate the models of the multiple properties of the target product. For each model, several data processing methods are empirically tested to optimize the accuracy and precision of the models.



4

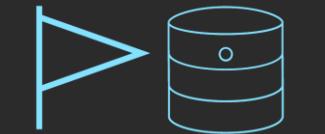


Application Development Roadmap



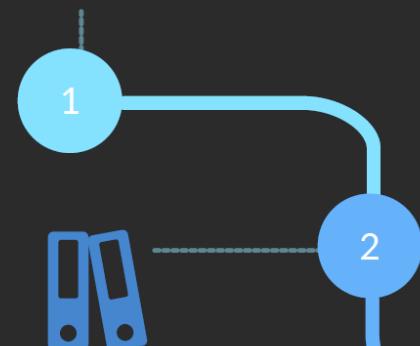
Pack into a predictive application

All the spectral models are packed into an application. All the technical and user-required documentation is generated along with the necessary licensing information.



Collect calibration data

Legacy and/or new data to develop calibrations is collected.



Organize data

Cleaning/filtering, sorting, harmonizing, pooling, storing and document all the calibration data gathered.

1

2



Calibrate spectral models

Use the calibration data to calibrate the models of the multiple properties of the target product. For each model, several data processing methods are empirically tested to optimize the accuracy and precision of the models.

4

5



Validation test

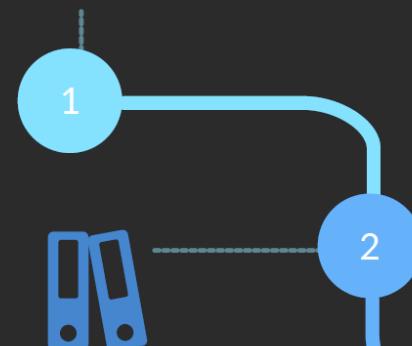
A specialist (different from the application developer), test that the application can be used. In this step the application is also tested with real samples to verify if reasonable results returned.

Application Development Roadmap



Collect calibration data

Legacy and/or new data to develop calibrations is collected.



Organize data

Cleaning/filtering, sorting, harmonizing, pooling, storing and document all the calibration data gathered.



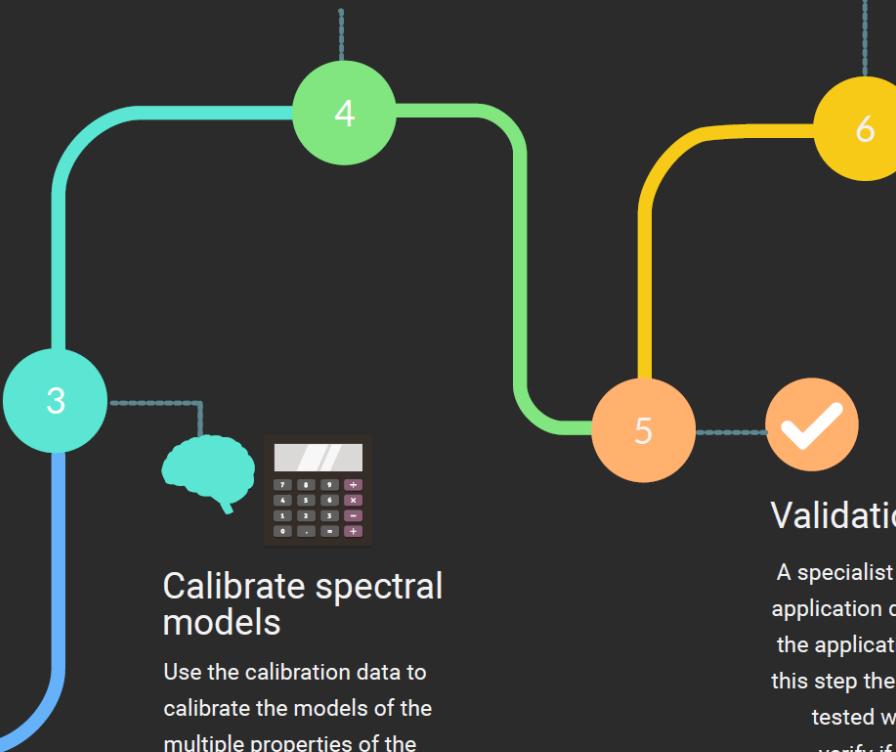
Pack into a predictive application

All the spectral models are packed into an application. All the technical and user-required documentation is generated along with the necessary licensing information.



Application release

The application is made available along with its technical documentation. Its availability is communicated to the stakeholders.



Calibrate spectral models

Use the calibration data to calibrate the models of the multiple properties of the target product. For each model, several data processing methods are empirically tested to optimize the accuracy and precision of the models.



Validation test

A specialist (different from the application developer), test that the application can be used. In this step the application is also tested with real samples to verify if reasonable results returned.



Application Development Roadmap



Collect calibration data

Legacy and/or new data to develop calibrations is collected.



Organize data

Cleaning/filtering, sorting, harmonizing, pooling, storing and document all the calibration data gathered.



Challenges

- It is **expensive**.
- Labs should focus on **routine analyses** (not on learning and dealing with chemometrics).
- For **end users/labs**, many time-consuming iterations to get useful applications creates **frustration**.
- **Chemometrics** is a complex subject which requires **time and effort** to be learned/taught.

Application Development Roadmap



Collect calibration data

Legacy and/or new data to develop calibrations is collected.

1



Organize data

cleaning/filtering, sorting, harmonizing, pooling, storing and document all the calibration data gathered.



Pack into an predictive application

All the spectral models are packed into an application. All the technical and user-required documentation is generated along with the necessary licensing information.



Application release

The application is made available along with its technical documentation. Its availability is communicated to the stakeholders.

4



Calibrate spectral models

Use the calibration data to calibrate the models of the multiple properties of the target product. For each model, several data processing methods are empirically tested to optimize the accuracy and precision of the models.

6



Validation test

A specialist (different from the application developer), test that the application can be used. In this step the application is also tested with real samples to verify if reasonable results returned.

5

Challenges

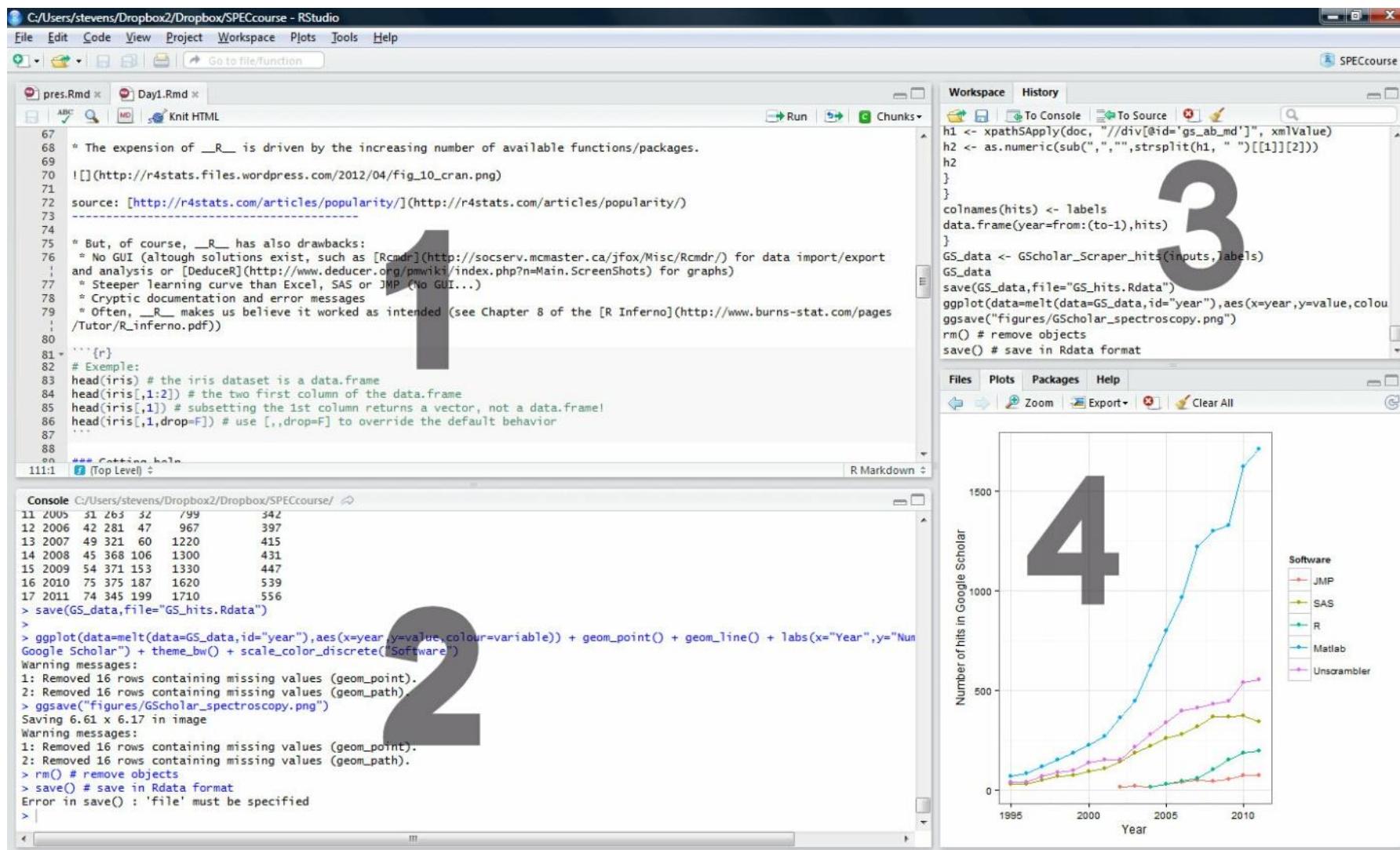
- It is **expensive to calibrate**, although it is extremely cost-efficient if it is released for routine analysis.
- Labs should be focused on **routine analyses** (not on learning and dealing with chemometrics)
- For **end users/labs**, many time-consuming iterations to get useful applications create **frustration**
- **Chemometrics** is a complex subject which requires time and effort to be learned/taught

Exercise

- Introduction to R
- Read in spectral data and visualization (1st script)

How to name your samples?

- Material name:
 - Core_[NAMEOFGROUP]
- Sample names:
 - [REPLICATE NUMBER]_[DEPTH increment]
 - Ex.: 01_0010
 - 01_1020
 - 02_1020



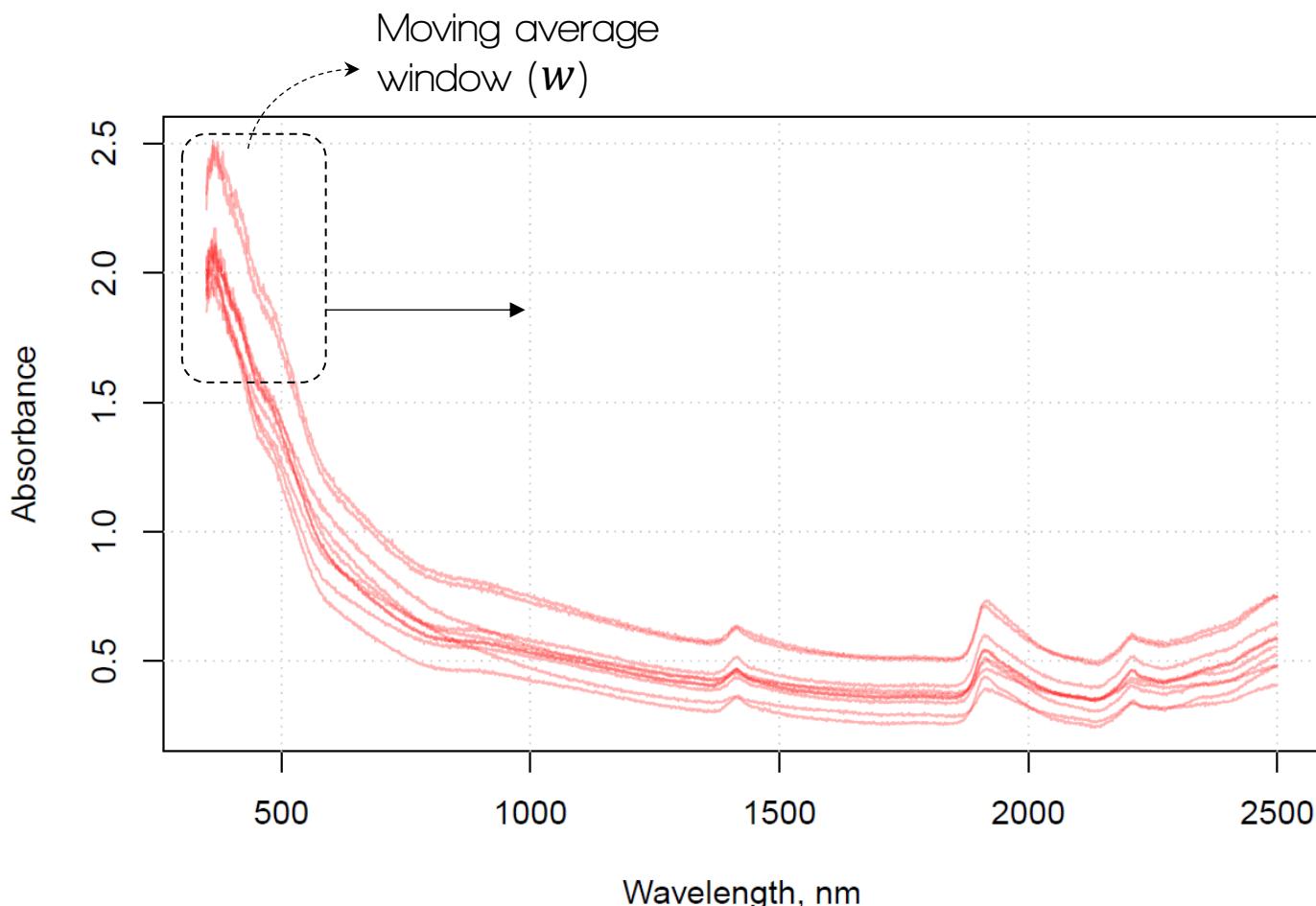
Signal processing aims to improve quality of the spectra before modeling (e.g. by removing undesired spectral variation). In this respect, it improves also the results of quantitative and qualitative spectral models.

The methods must be carefully chosen in order to avoid removing important physical information.

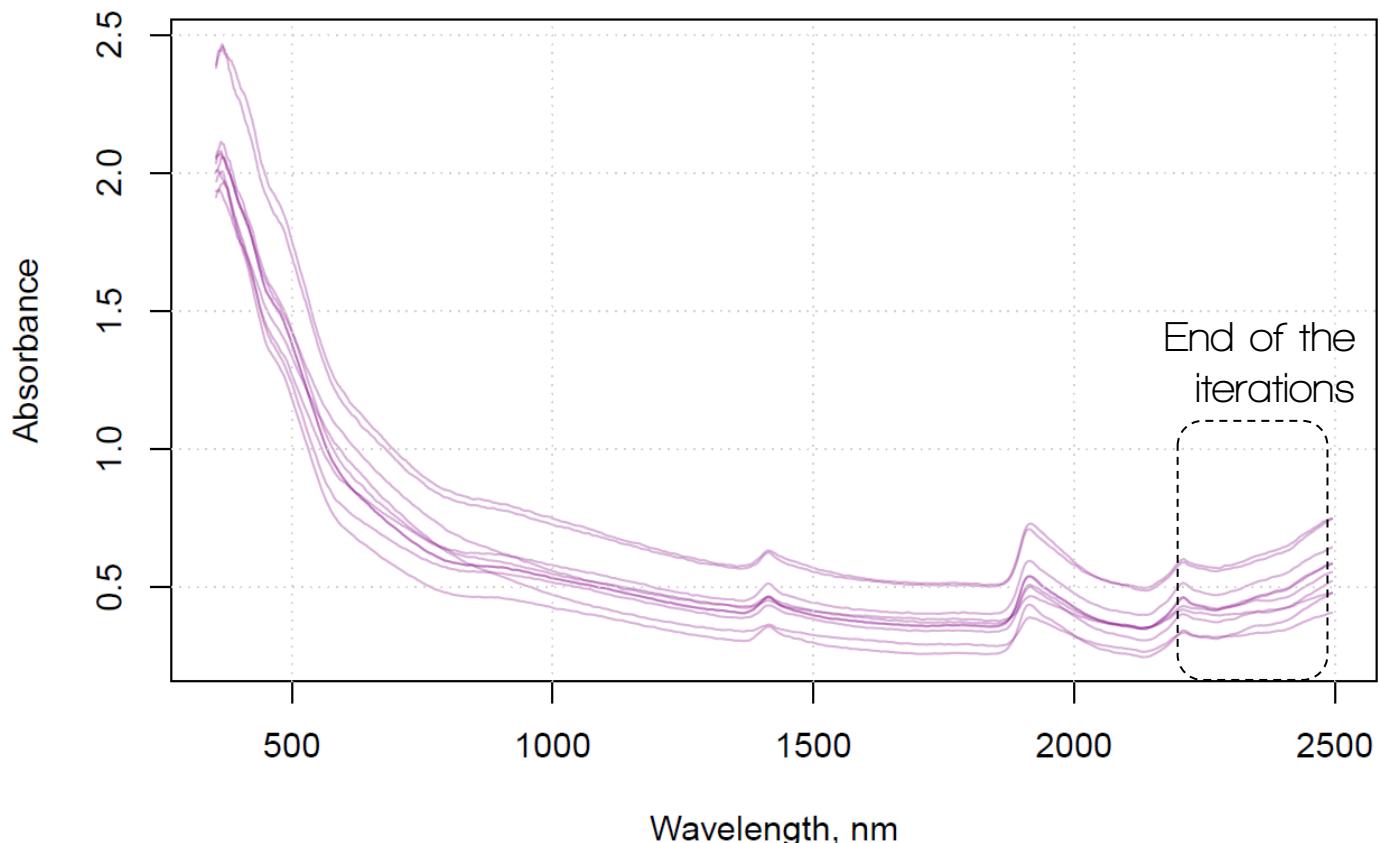
The most common signal processing methods are:

- Splice correction
- Spectra transformation (e.g. reflectance to absorbance)
- Noise removal
- Resampling
- Differentiation
- Scatter correction
- Centering and scaling
- Continuum removal

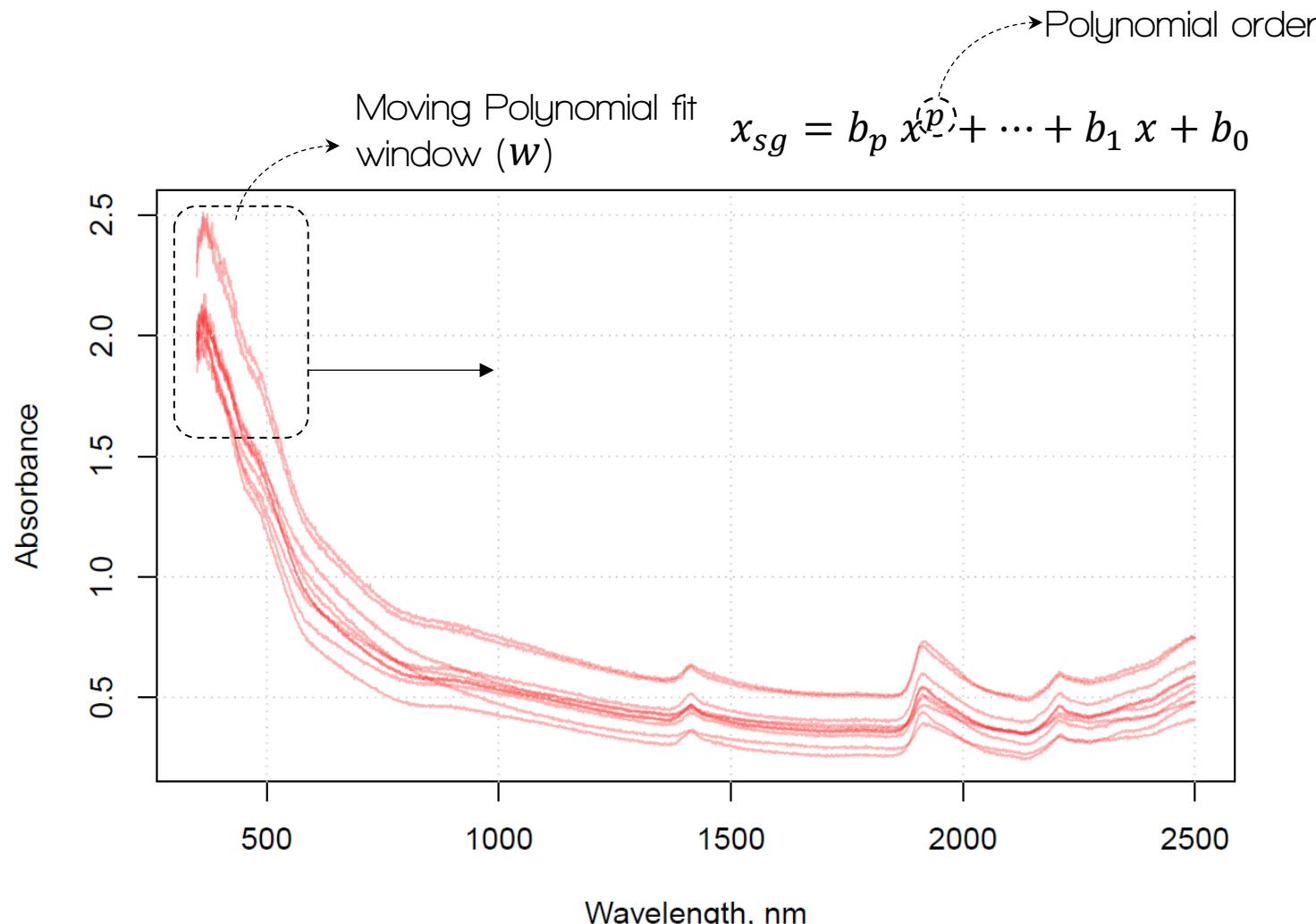
Moving average



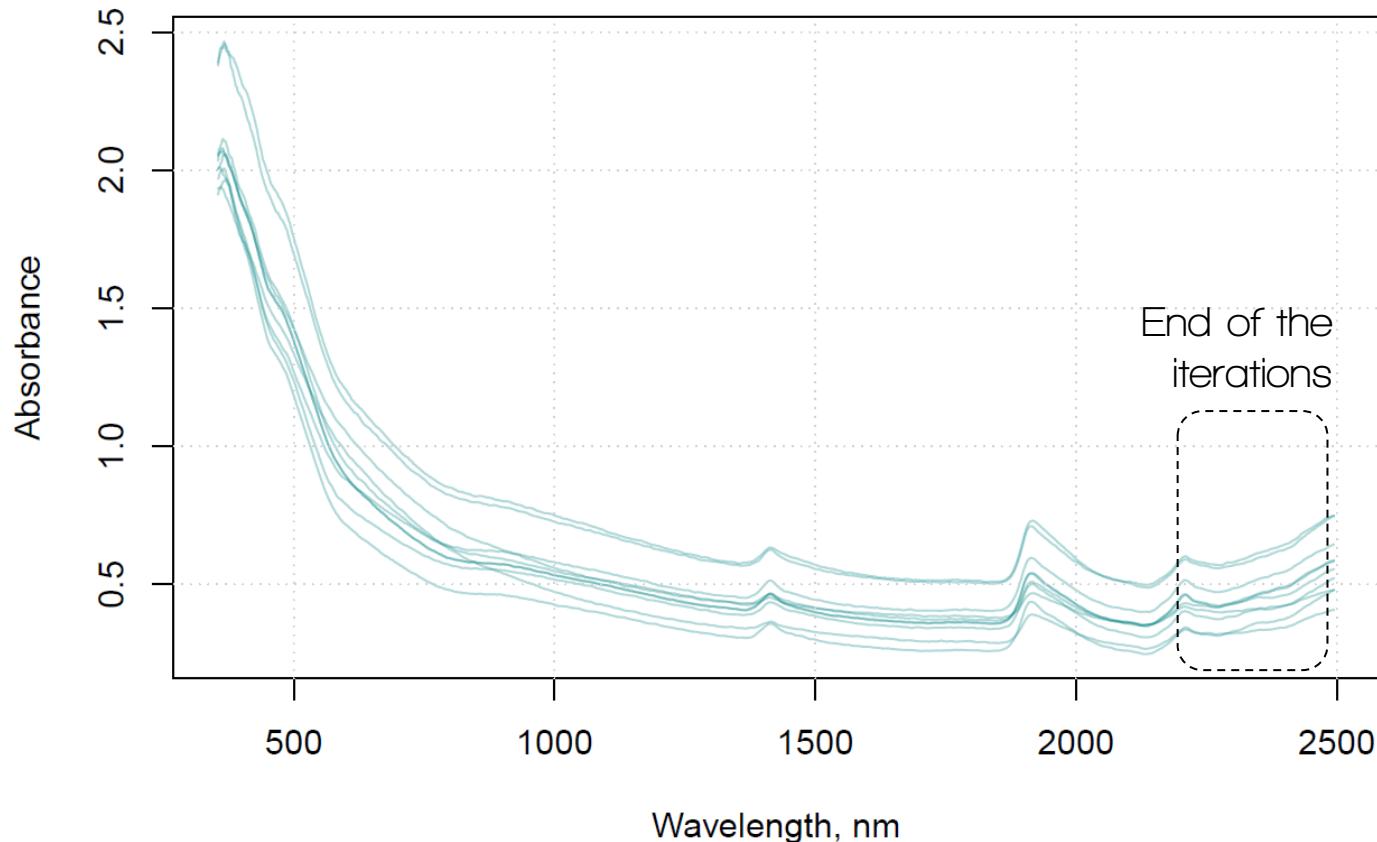
Moving average



Savitzky–Golay filter

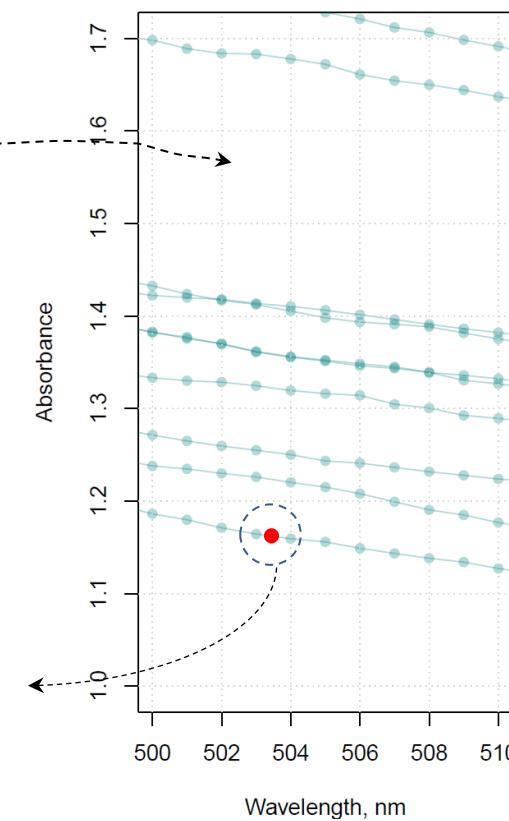
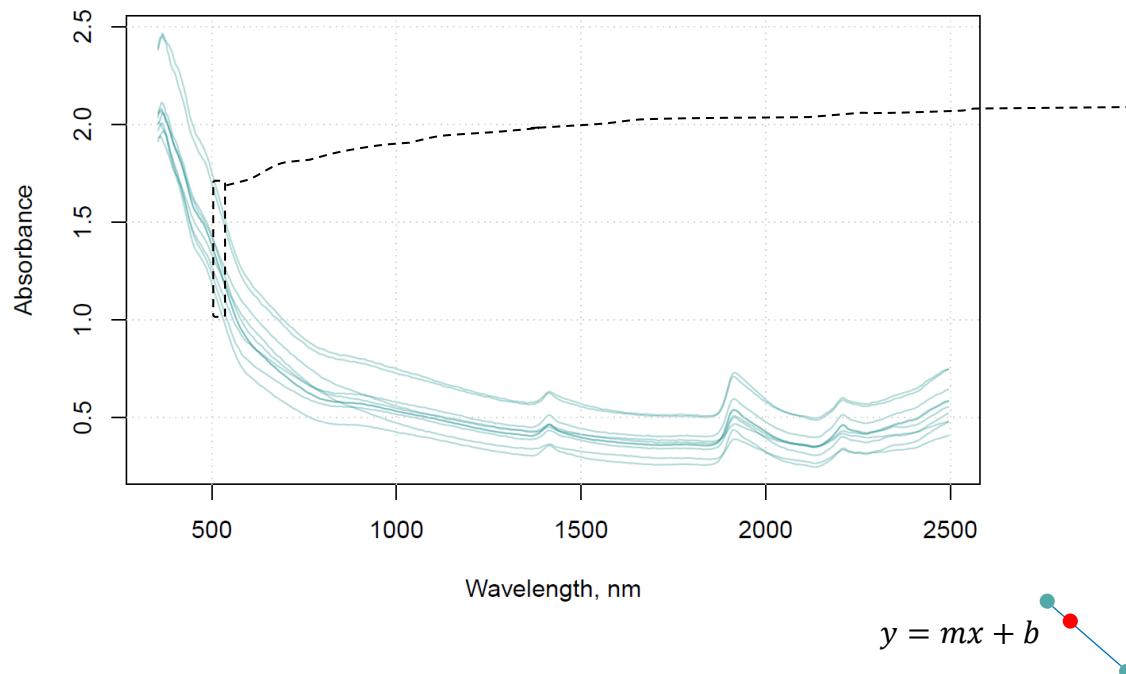


Savitzky–Golay filter



Resampling (also known as interpolation)

Example: if the first wavelength of the new set of wavelengths (at which the spectra needs to be resampled) is **504.3 nm**, the two wavelengths at the right and left sides of the new point are taken. Then a linear model is fitted between these two points and this model is then used to estimate the intensity at the new wavelength



Advantages

- Reduce of baseline offset
- Helps to resolve absorption overlapping
- Compensates for instrumental drift
- Enhances small spectral absorption features
- Often increases the predictive accuracy for complex datasets

Drawbacks

- Risk of overfitting
- Amplifies noise, smoothing required
- Less robust model and increase uncertainty in model coefficients
- Complicate spectral interpretation
- Remove the baseline! (albedo variation)

First and second derivatives of a spectrum can be computed with the finite difference method (difference between two subsequent data points), **assuming that the band width is constant**:

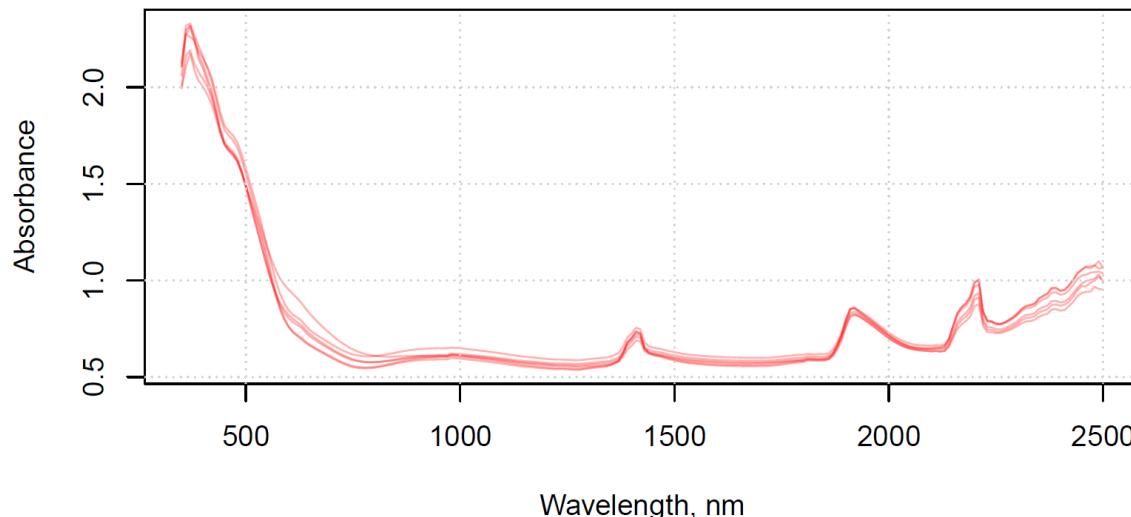
$$\text{First derivative } x'_i = x_i - x_{i-1}$$

$$\text{Second derivative } x''_i = x_{i-1} - 2x_i + x_{i+1}$$

Signal processing

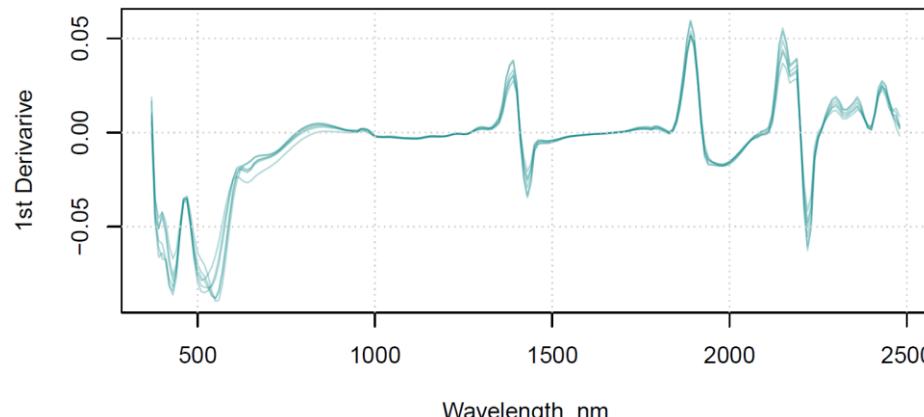
Differentiation

Original absorbance spectra of a soil profile

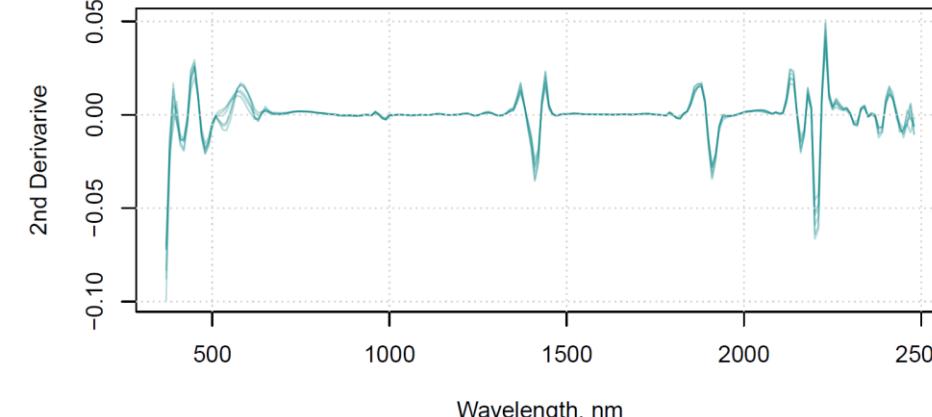


Using the **savitzkyGolay** function both the first and second derivative spectra, can be computed

First derivative spectra



Second derivative spectra



Multiplicative Scatter Correction (msc)

- It mitigates the effect of scattered light (which is multiplicative) on the NIR spectra
- **msc** is a method that assumes for a given set of spectra, the scattering and offset effects can be removed from each spectrum.
- The **msc** model is determined by regressing each spectrum onto the mean spectrum. The model can be formally written as:

$$x_i = a_i + b_i \bar{x} + \varepsilon_i$$

where ε_i represents the chemical information of the spectrum while a_i and b_i represent the offset and the scatter respectively. Both a_i and b_i are found by using the Least Squares method. The **msc** corrected spectrum can be then computed as:

$$\text{msc}(x_i) = \frac{x_i - a_i}{b_i}$$

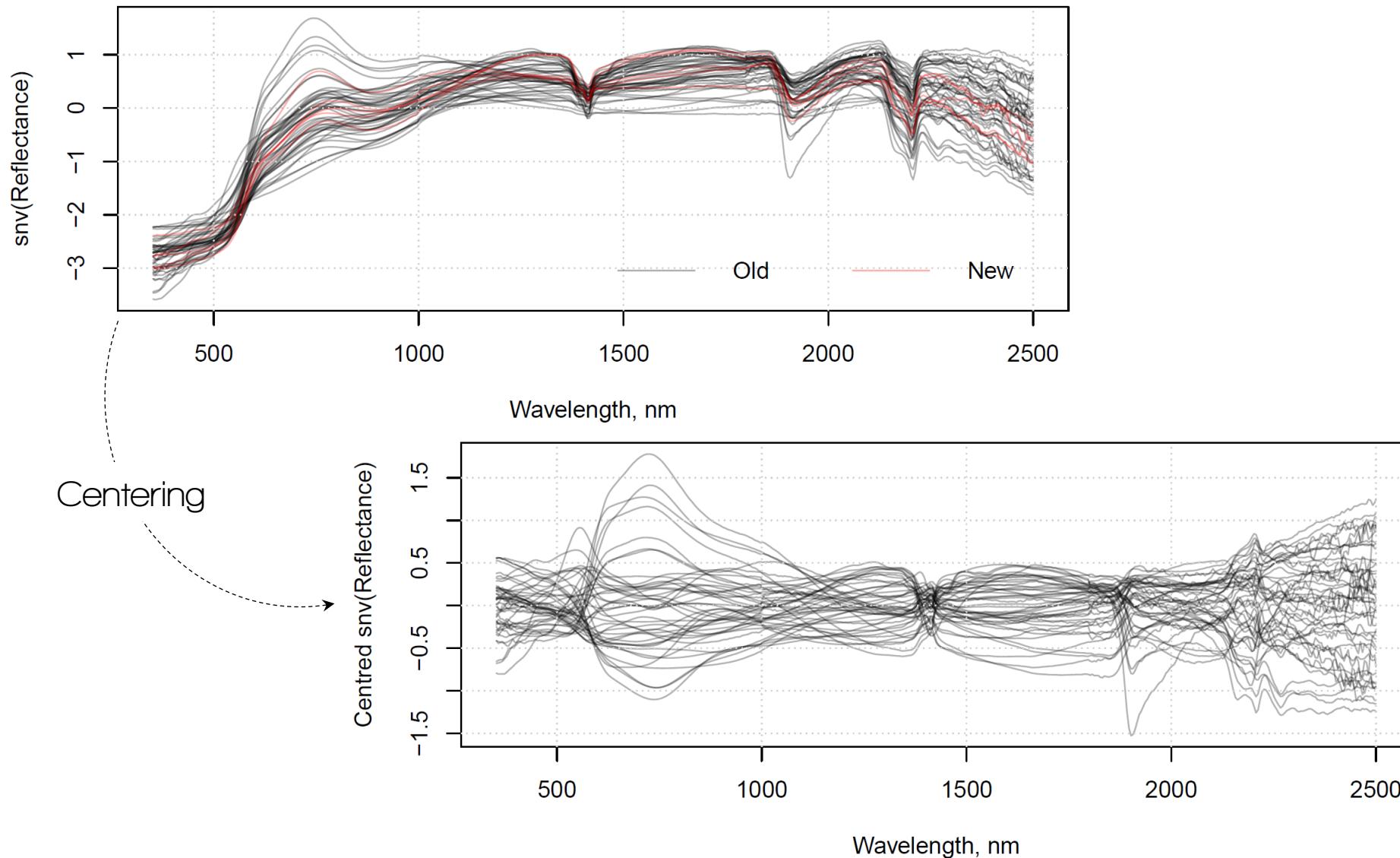
The **msc** function pf the **pls** package can be used to apply the Multiplicative Scatter Correction

- *Centering*: a given matrix is transform into a matrix with columns that have zero mean. This is achieved by subtracting the mean vector of the matrix columns from each row of the original matrix. From a statistical point of view, *centering* eliminates the need for an intercept in the model and makes the model more simple and interpretable. This operation is **strongly** advised for most regression approaches, although it does not imply that in all cases the model performance will improve
- *Scaling*: a given matrix is transform into a matrix with columns that have the same variance. This is achieved by dividing each row of the original matrix by the vector of standard deviations of the matrix columns. The purpose is to get all the variables . . . on the same scale.
- Both centering and scaling are usually applied after data transformation.

The `scale` function can be used for applying either centering and scaling or both

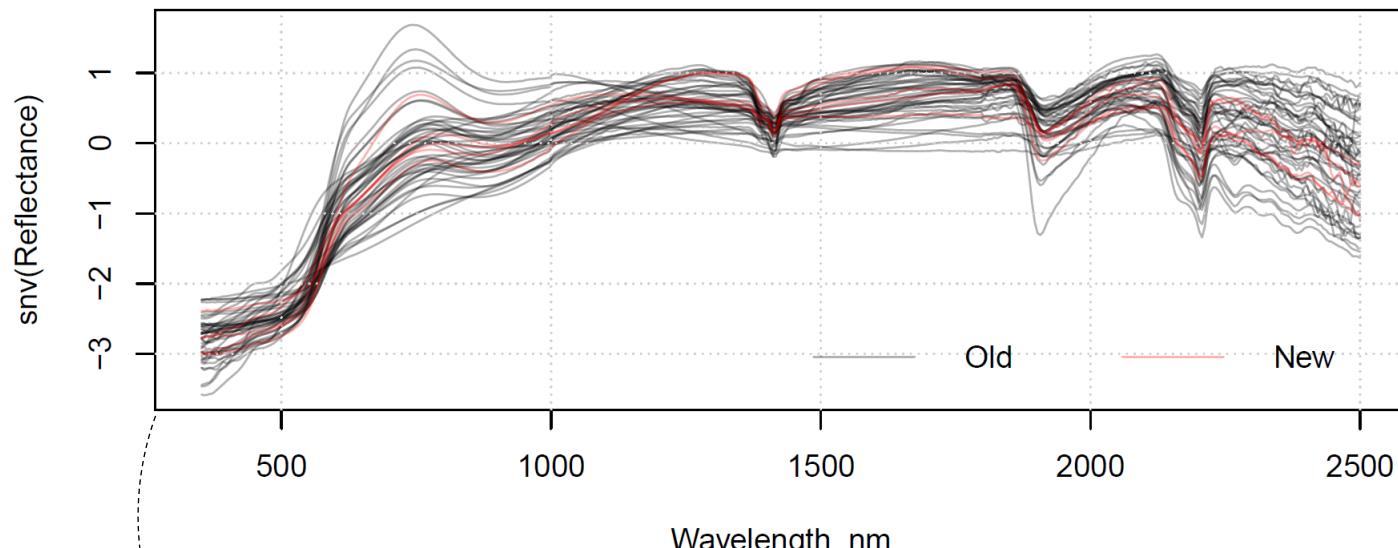
Signal processing

Centering and scaling



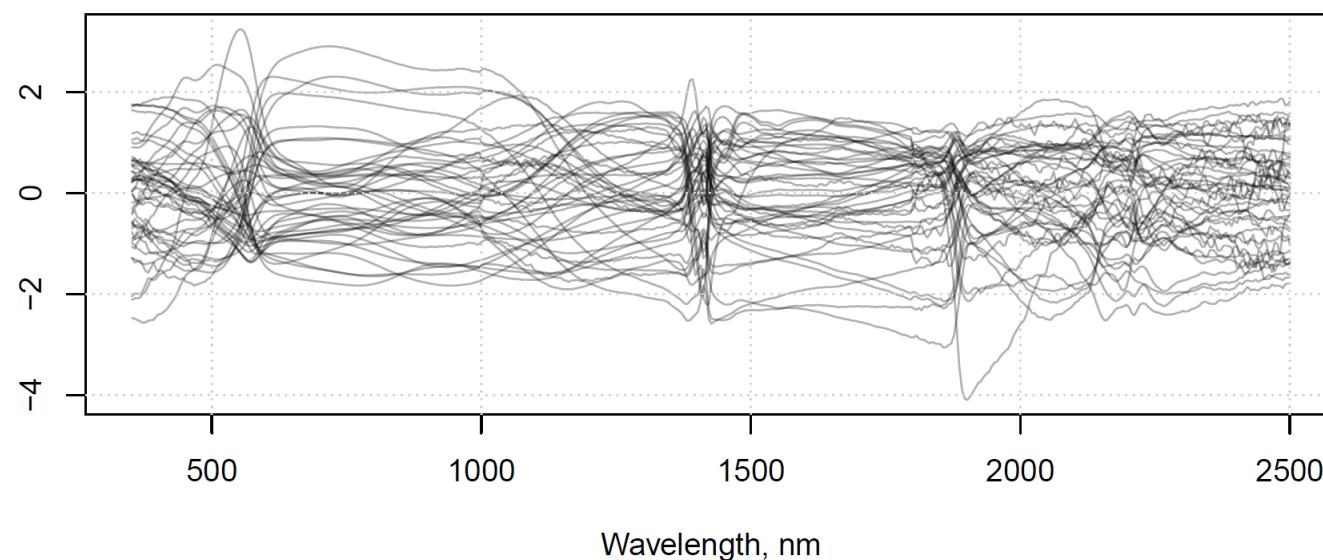
Signal processing

Centering and scaling



Centering
and
scaling

→ Centred and scaled snv(Reflectance)



To *scale* or not to *scale*?

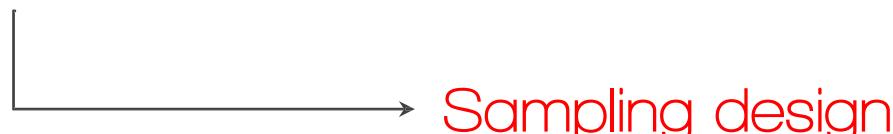
- Scaling will down-scale variables with large variation, but can inflate exaggeratedly variables with small variance (e.g. variables with no variation but small noise).
- Decision depends on the data pre-treatment also: is it meaningful?
- For instance, variation in some variables is greater than in other for good reasons (e.g. absorption features in continuum removed spectra, 1st derivative), so that *scaling* is not appropriate and will reduce model interpretability
- *Centering and scaling* is sensitive to outliers. A robust scaling can be computed using the median instead of the mean and the median absolute deviation (MAD) instead of the standard deviation

How many observations (samples) should be included in a calibration set in order **to efficiently develop the best NIR models?**

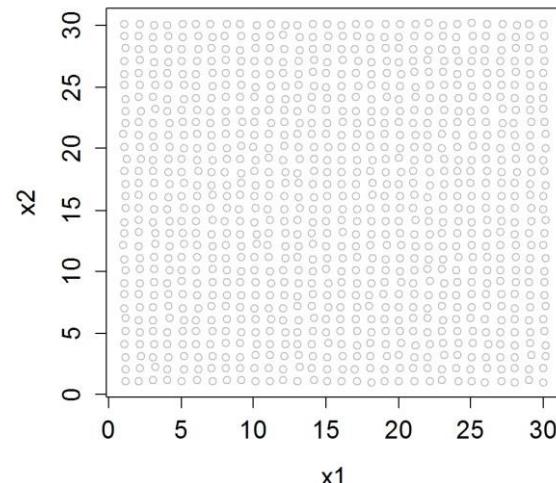
How many samples should be used for calibration?



How to select the calibration samples?



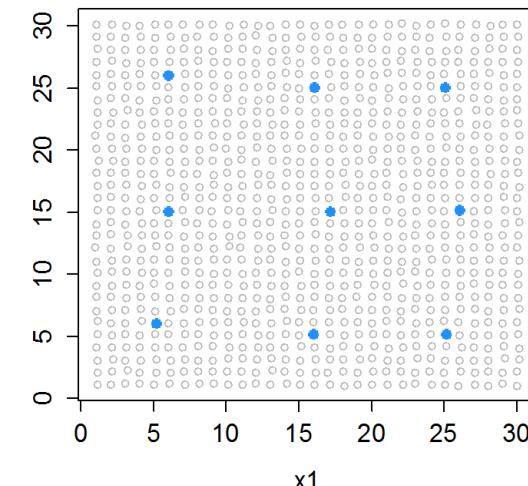
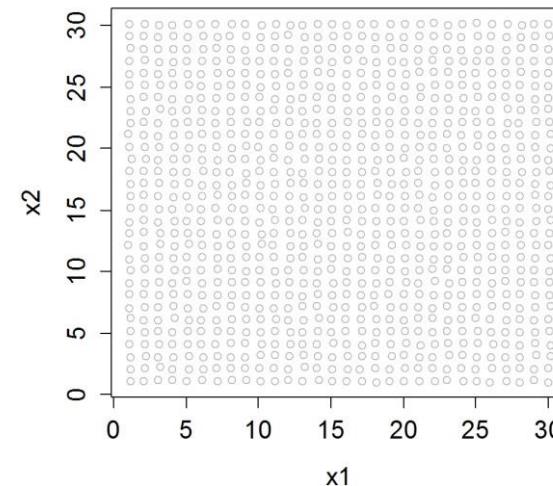
Sample representativeness



Calibration sampling

Sampling design

Sample representativeness



Kennard-Stone sampling (KSS)

The KSS algorithm (Kennard and Stone, 1969) to select a (calibration) subset of m samples (S_{tr}) from a given set of N samples (S) works as follows:

1. Compute the distance matrix between all the samples in X
2. Find in S the two samples s_{tr1} and s_{tr2} , which are farthest apart, allocate them in S_{tr} and remove them from S .
3. Find in S the sample s_{tr3} , which is the most dissimilar to the ones already allocated in S_{tr} . Allocate s_{tr3} in S_{tr} and then remove it from S . Note that the dissimilarity between S_{tr} and each s_i is given by the minimum distance of any sample allocated in S_{tr} to each s_i .
4. Repeat the step $3m - 4$ times in order to select the remaining samples.

The KSS is usually applied to the PC scores of the spectral data

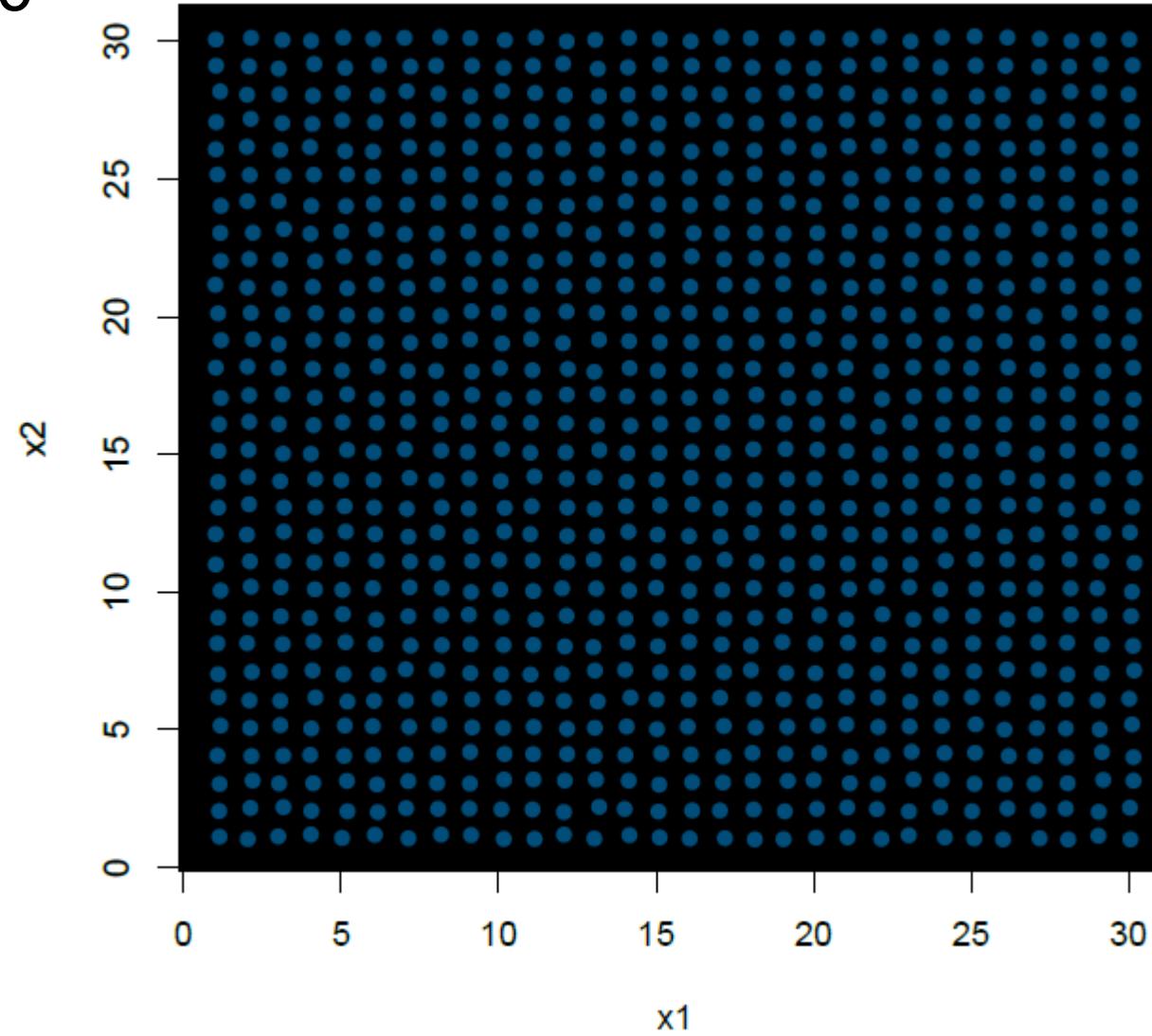
The **kenStone** function of the **prospectr** package can be used for KSS

Calibration sampling

Sampling design

$$n = \emptyset$$

Kennard-Stone sampling (KSS)

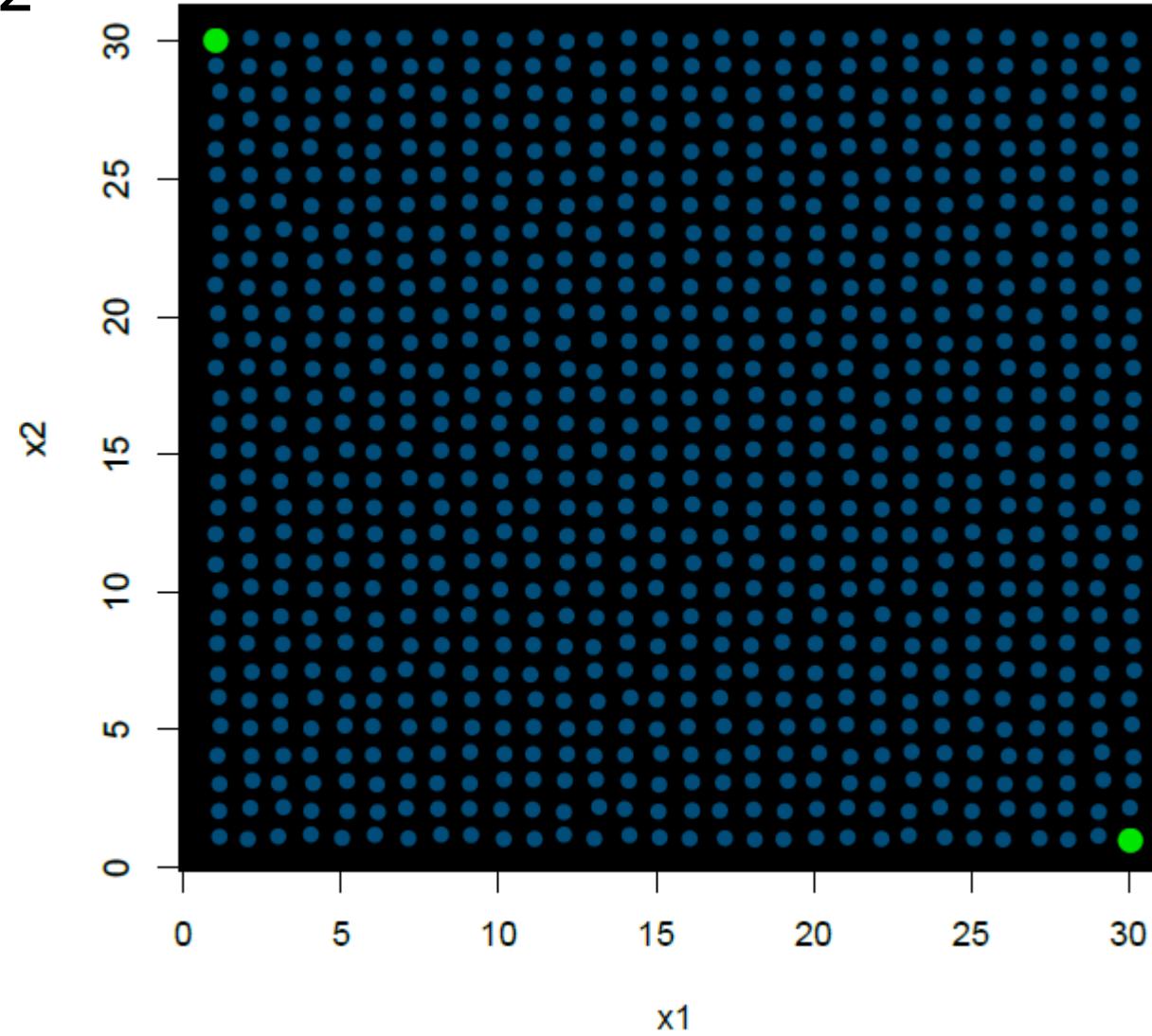


Calibration sampling

Sampling design

n = 2

Kennard-Stone sampling (KSS)

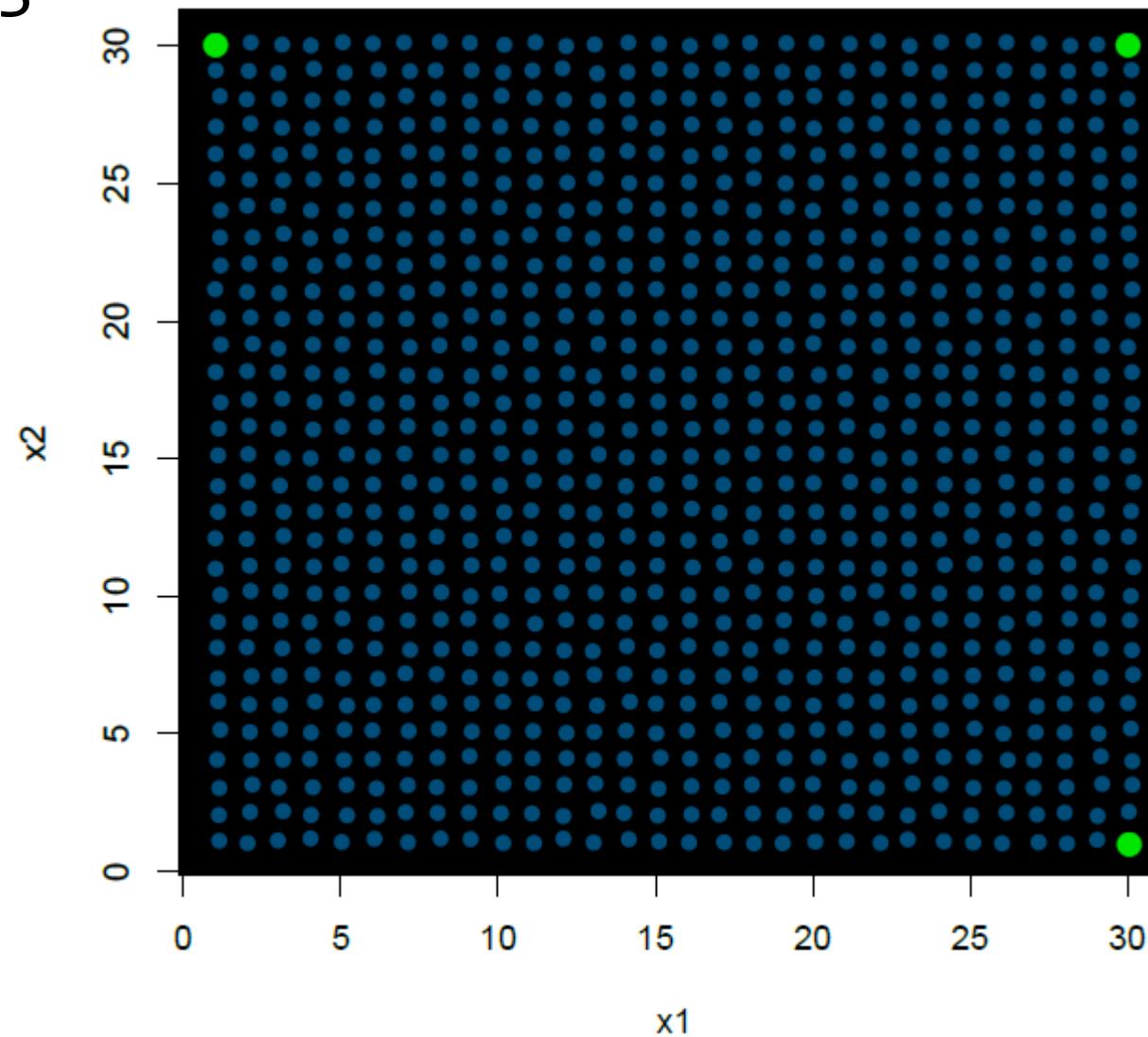


Calibration sampling

Sampling design

n = 3

Kennard-Stone sampling (KSS)

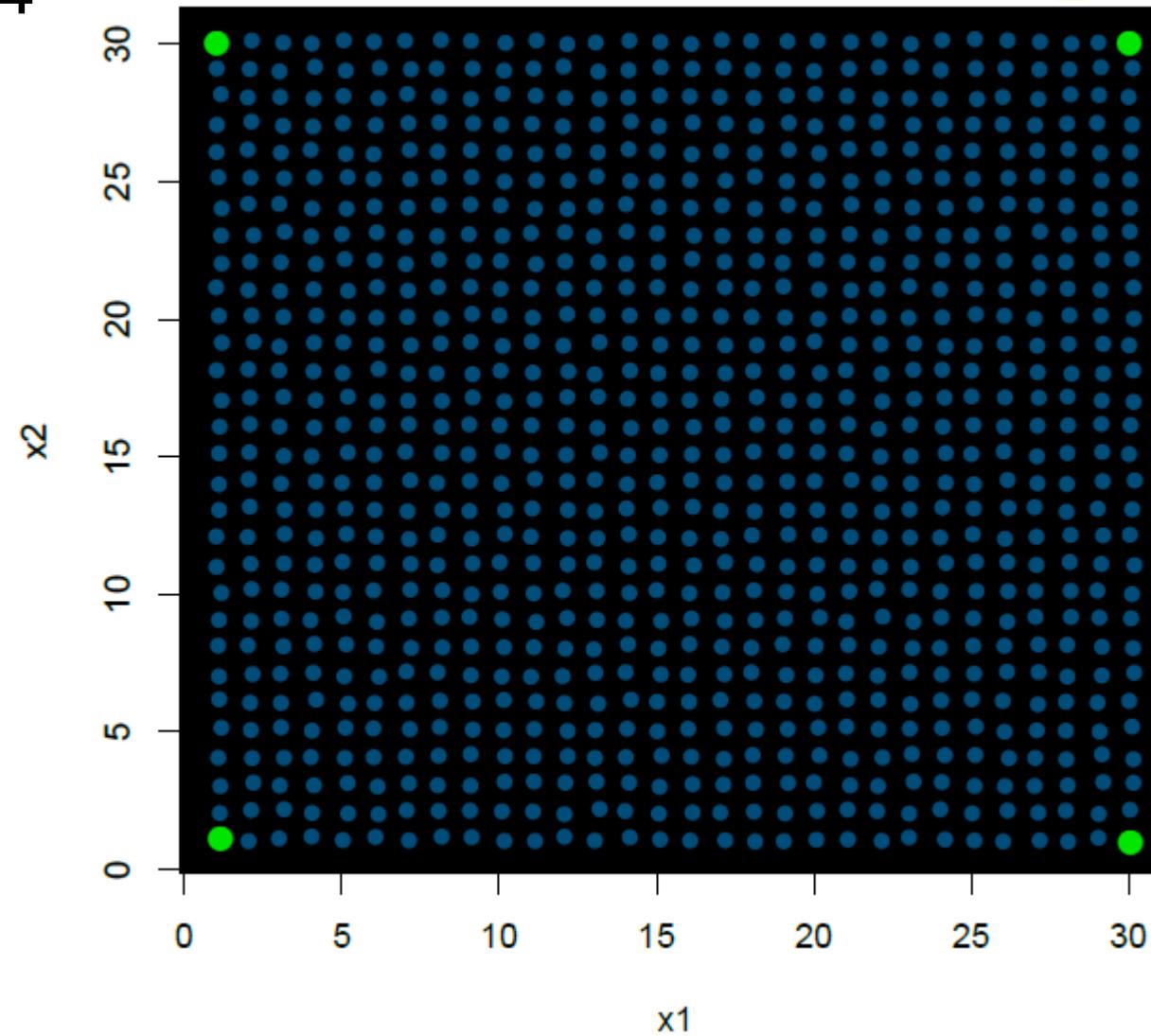


Calibration sampling

Sampling design

n = 4

Kennard-Stone sampling (KSS)

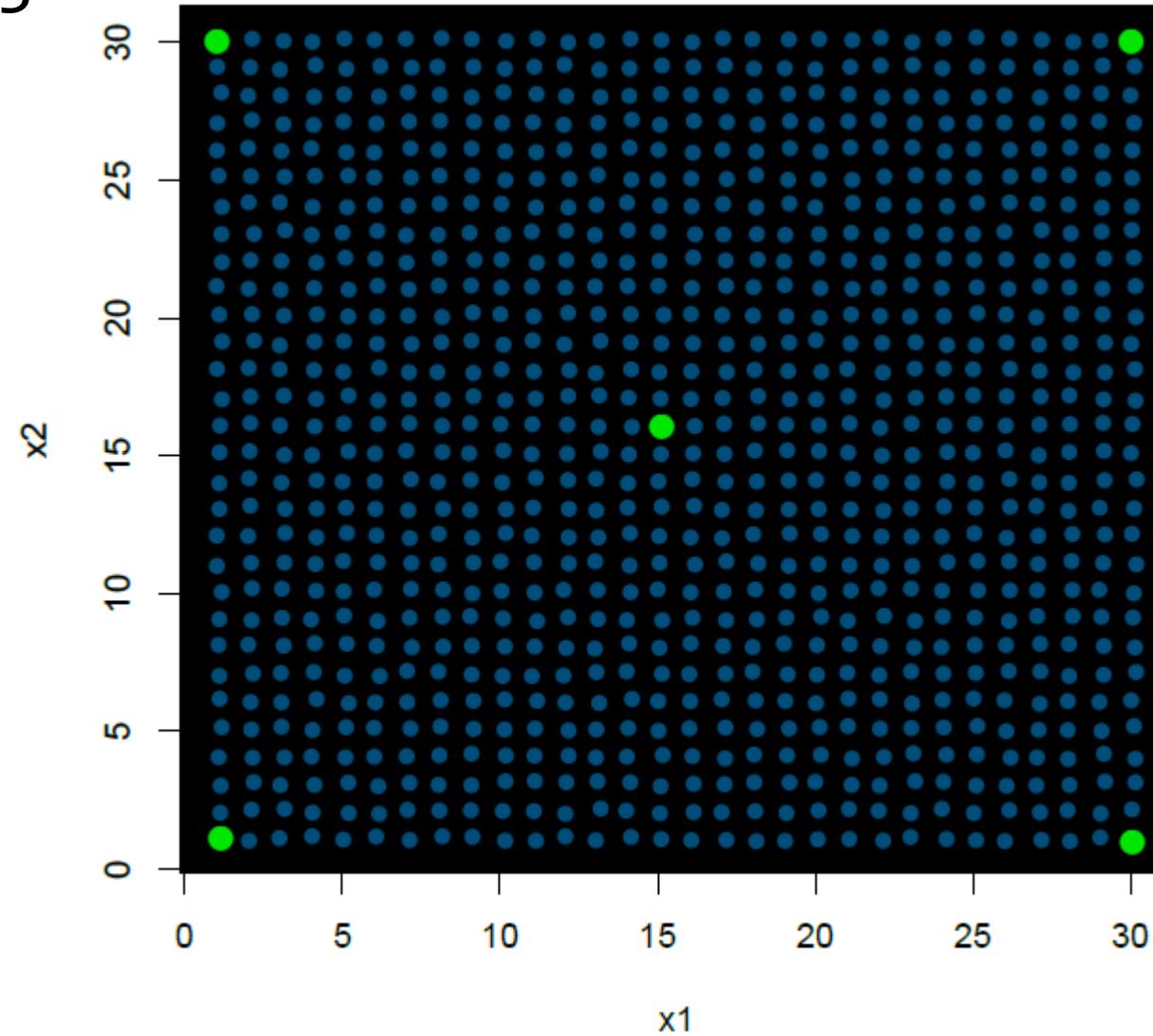


Calibration sampling

Sampling design

n = 5

Kennard-Stone sampling (KSS)

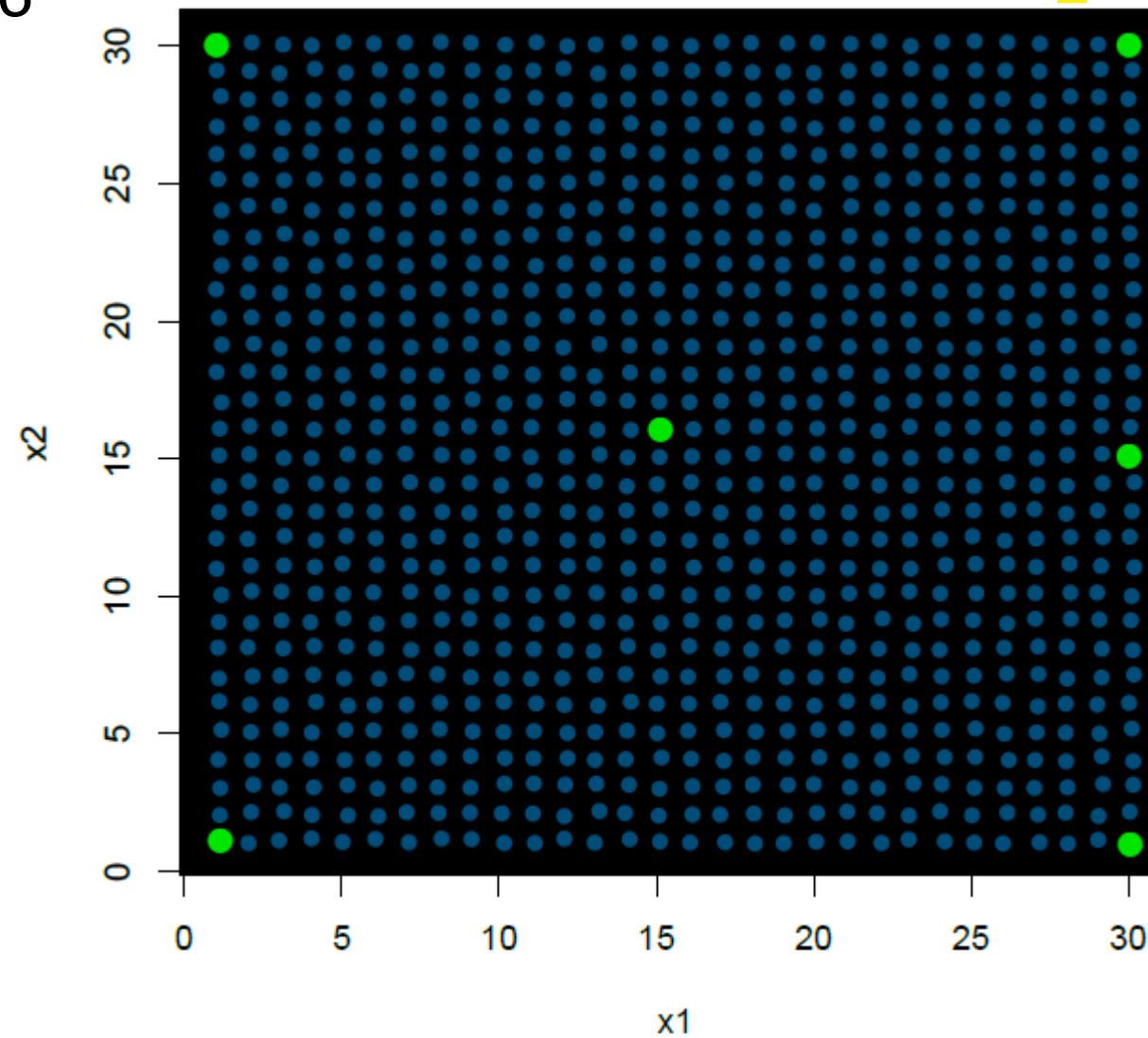


Calibration sampling

Sampling design

n = 6

Kennard-Stone sampling (KSS)

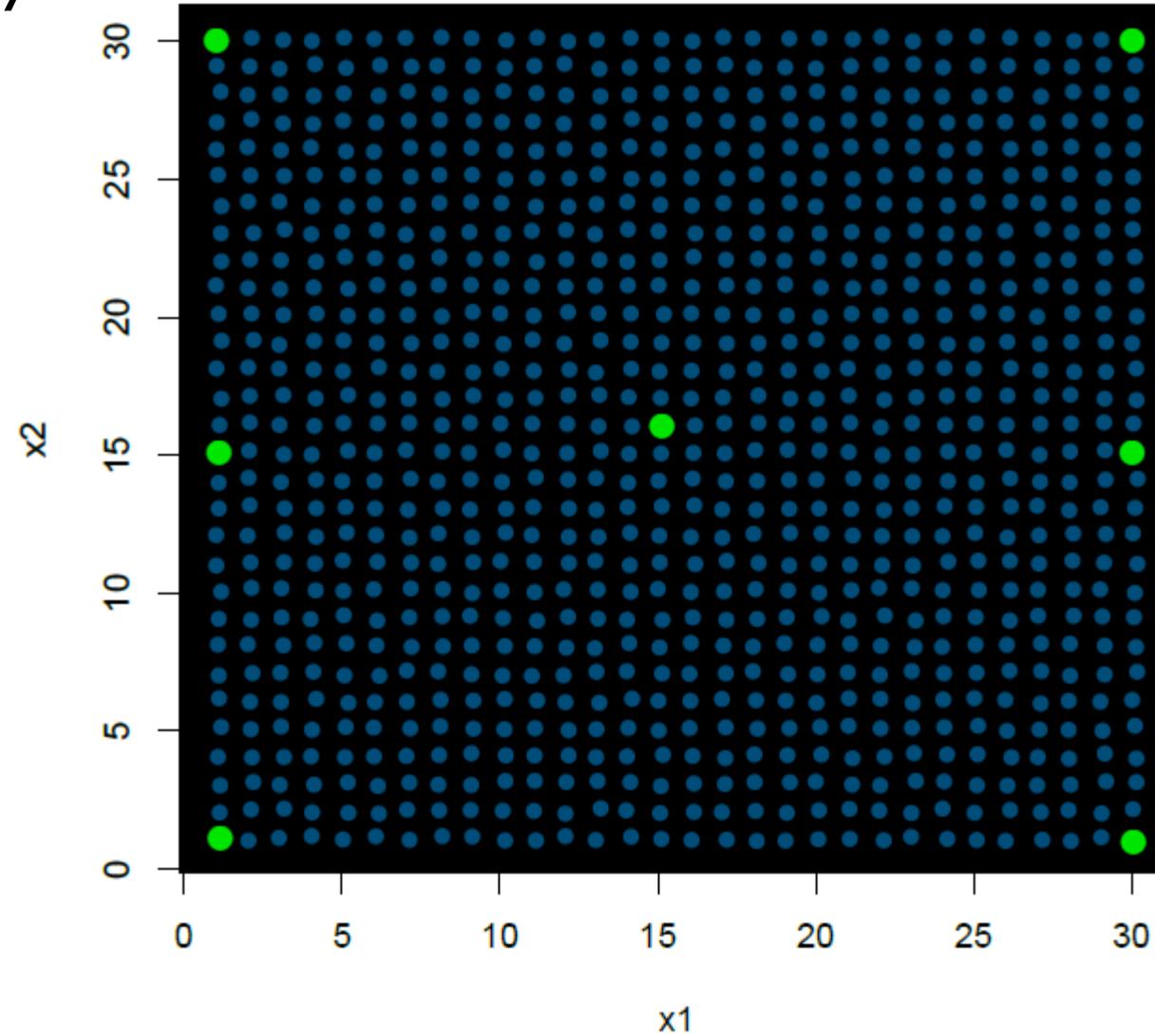


Calibration sampling

Sampling design

n = 7

Kennard-Stone sampling (KSS)

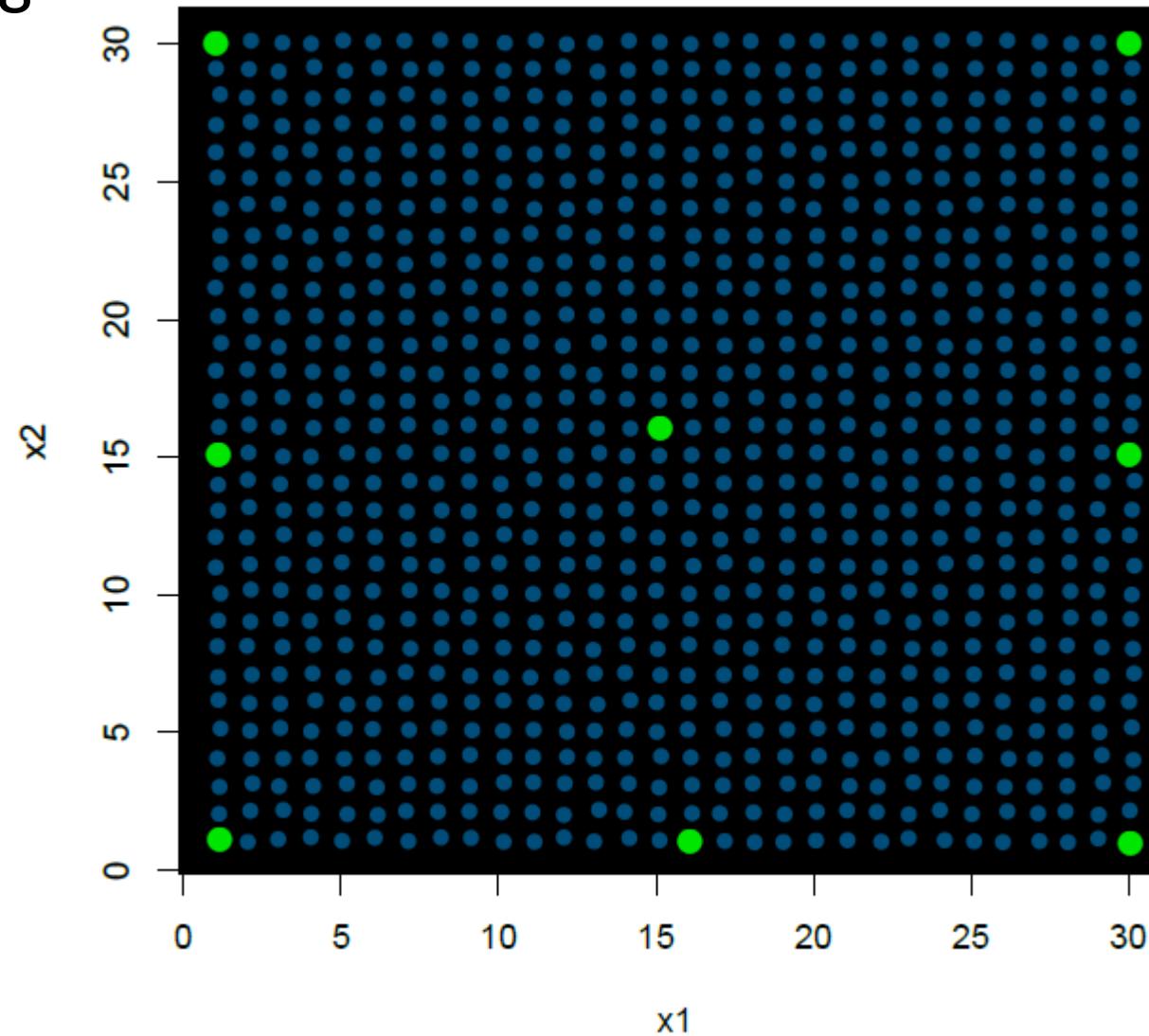


Calibration sampling

Sampling design

n = 8

Kennard-Stone sampling (KSS)

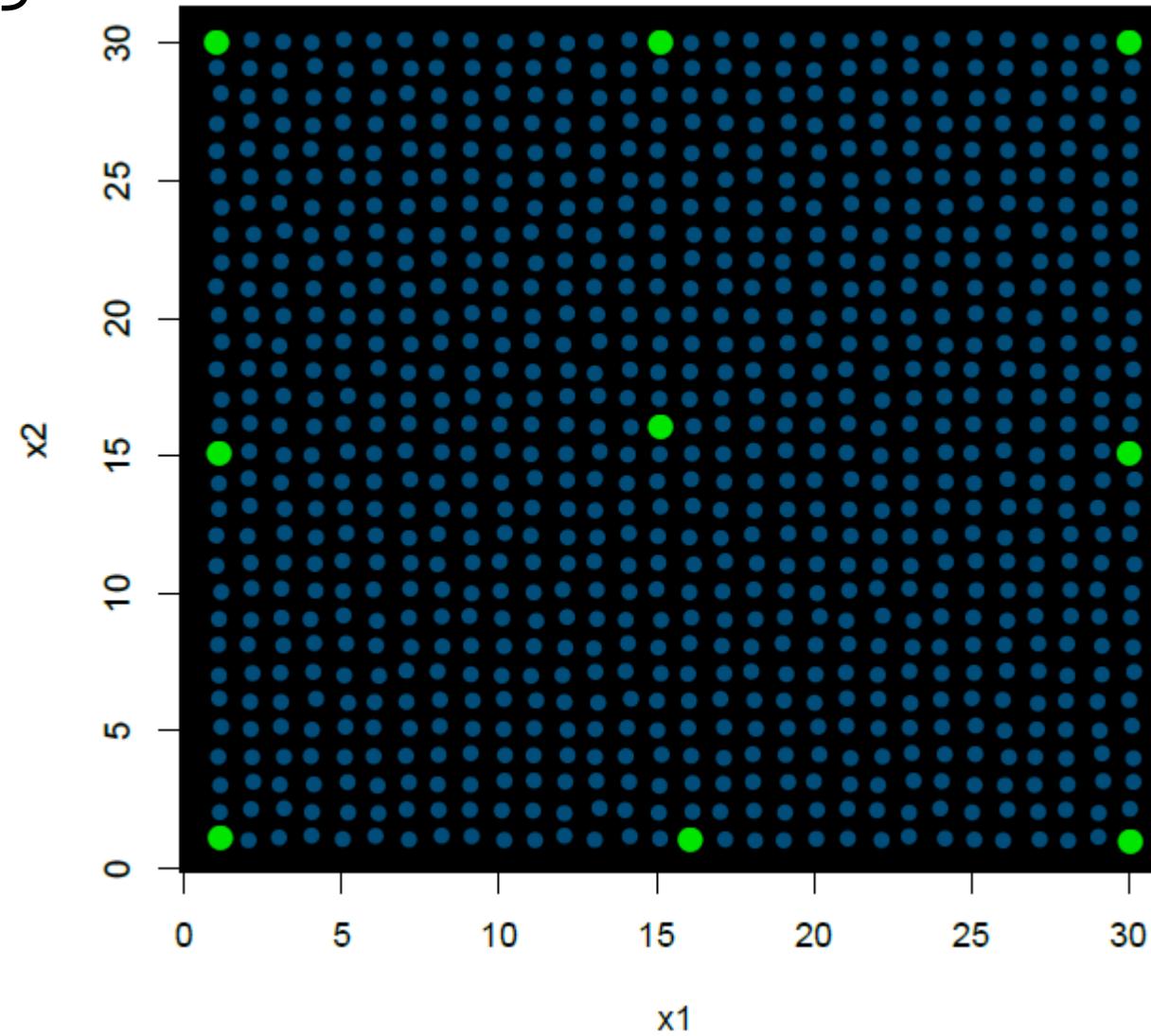


Calibration sampling

Sampling design

n = 9

Kennard-Stone sampling (KSS)

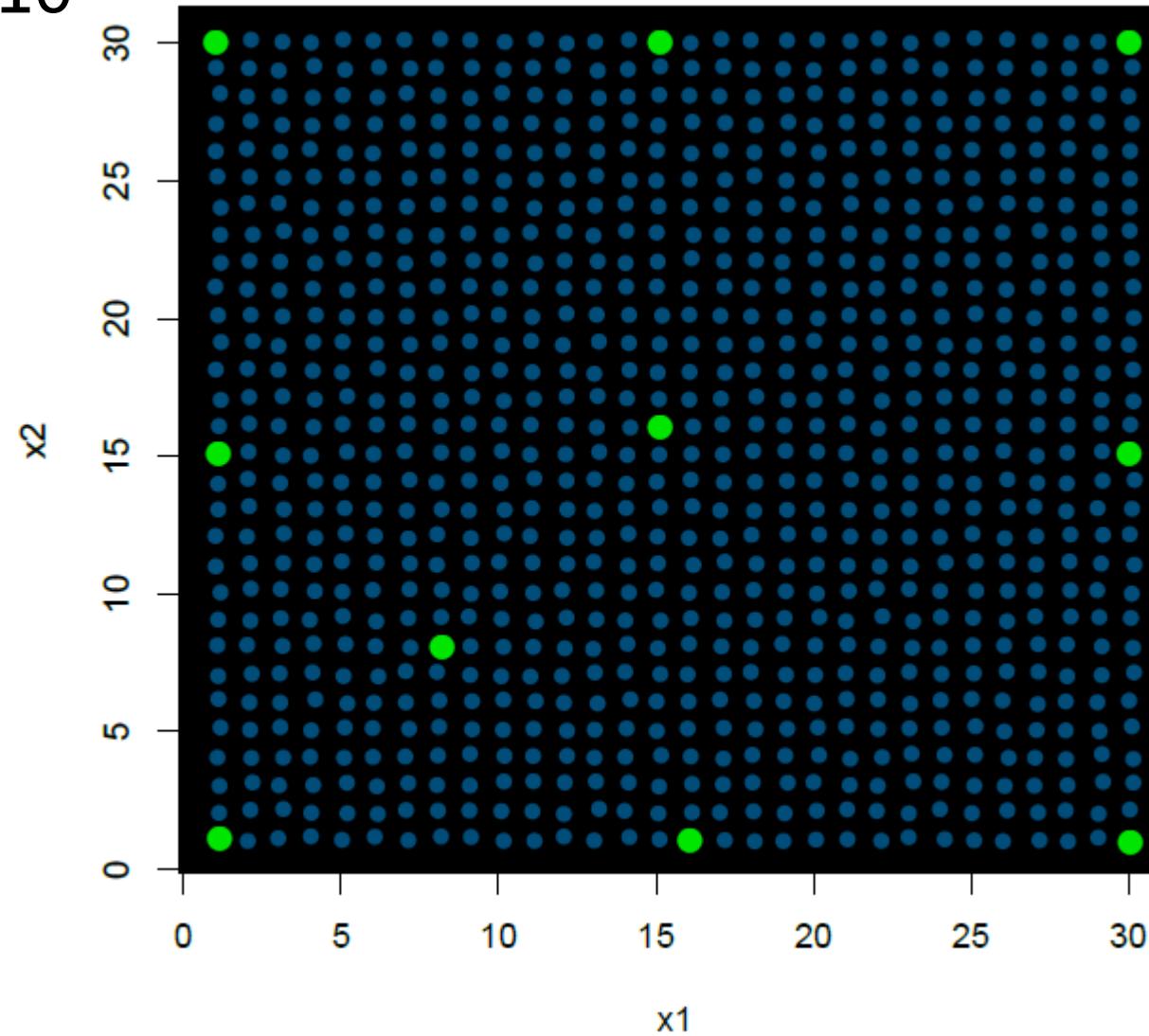


Calibration sampling

Sampling design

n = 10

Kennard-Stone sampling (KSS)

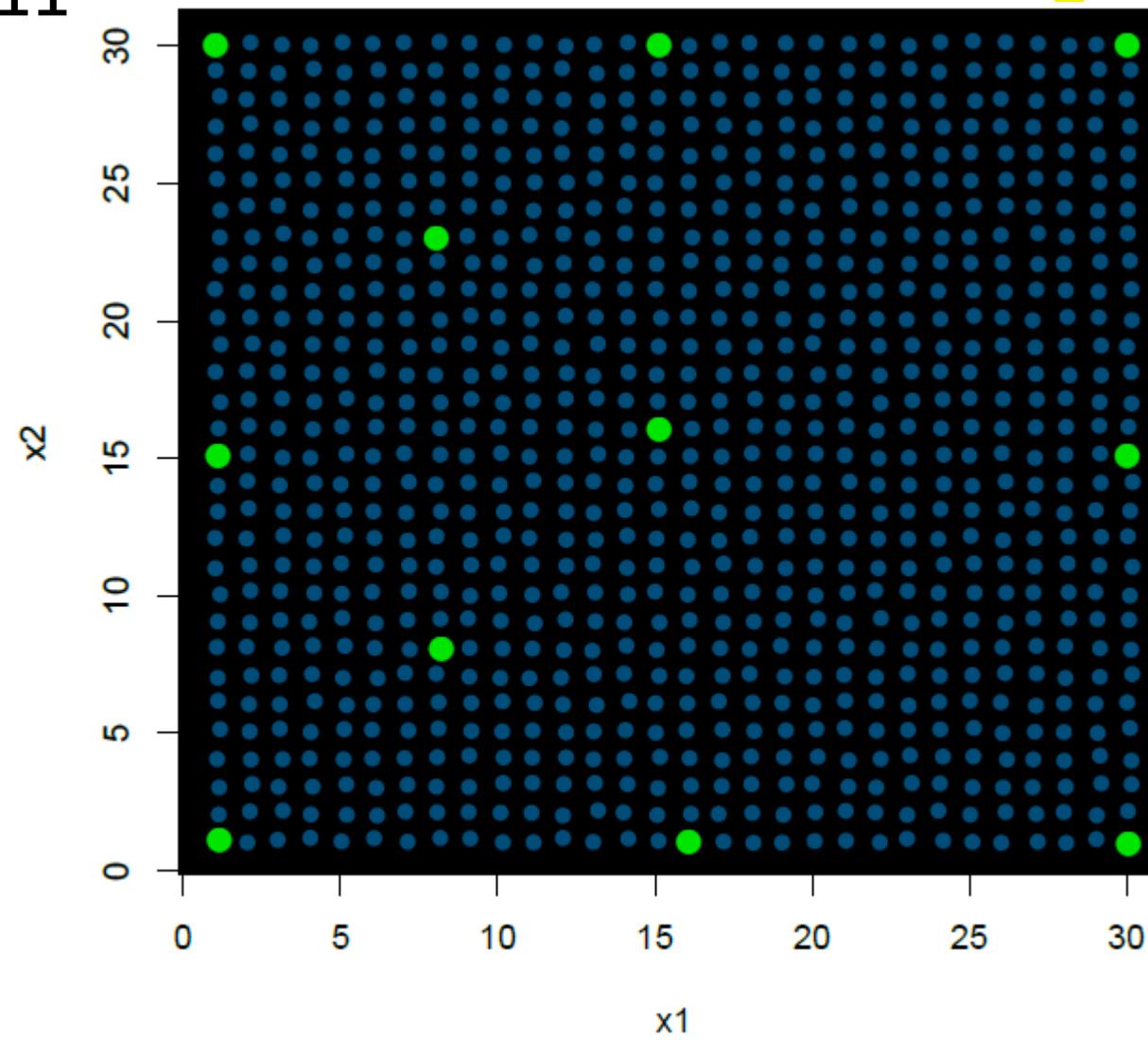


Calibration sampling

Sampling design

n = 11

Kennard-Stone sampling (KSS)

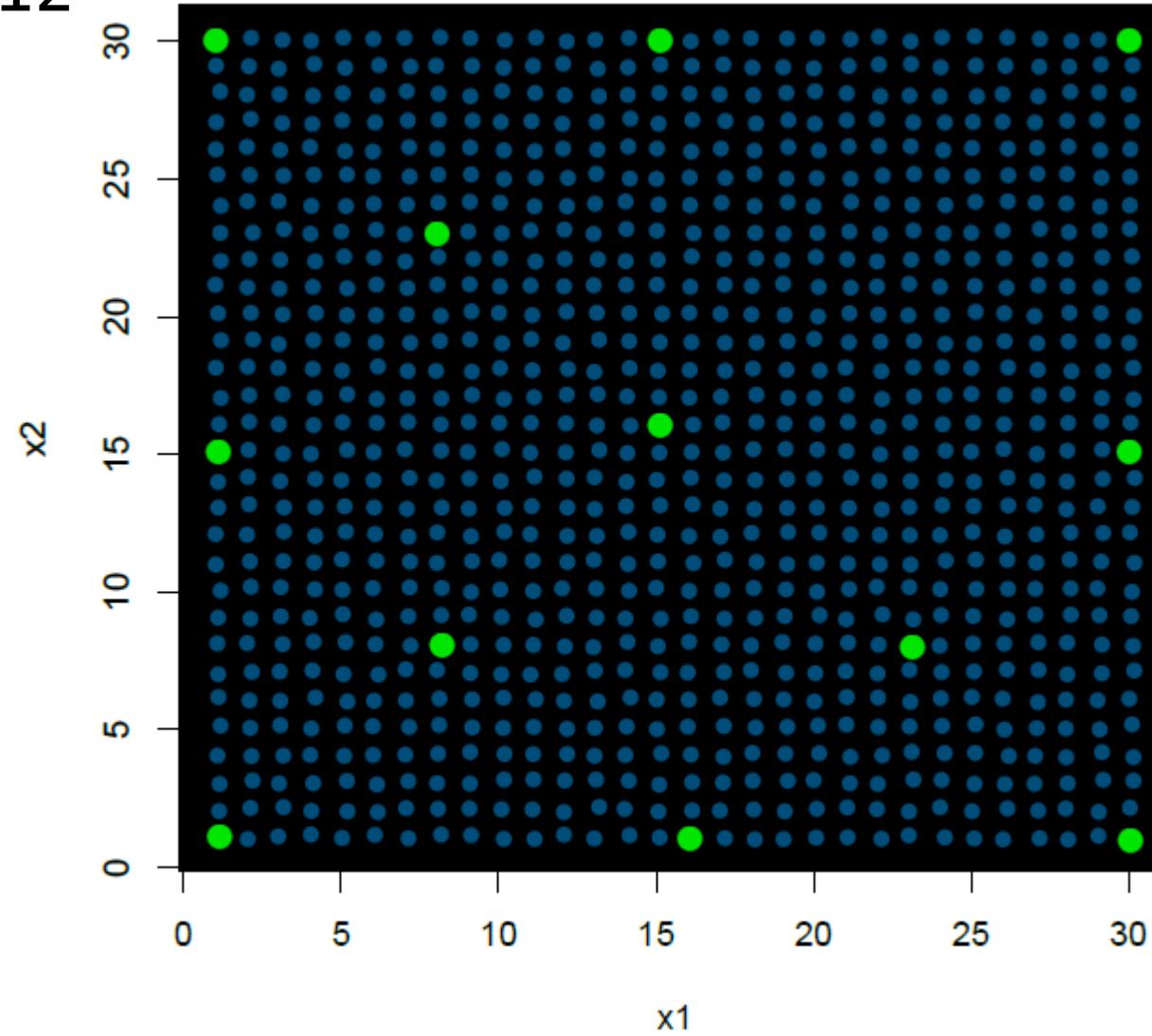


Calibration sampling

Sampling design

n = 12

Kennard-Stone sampling (KSS)

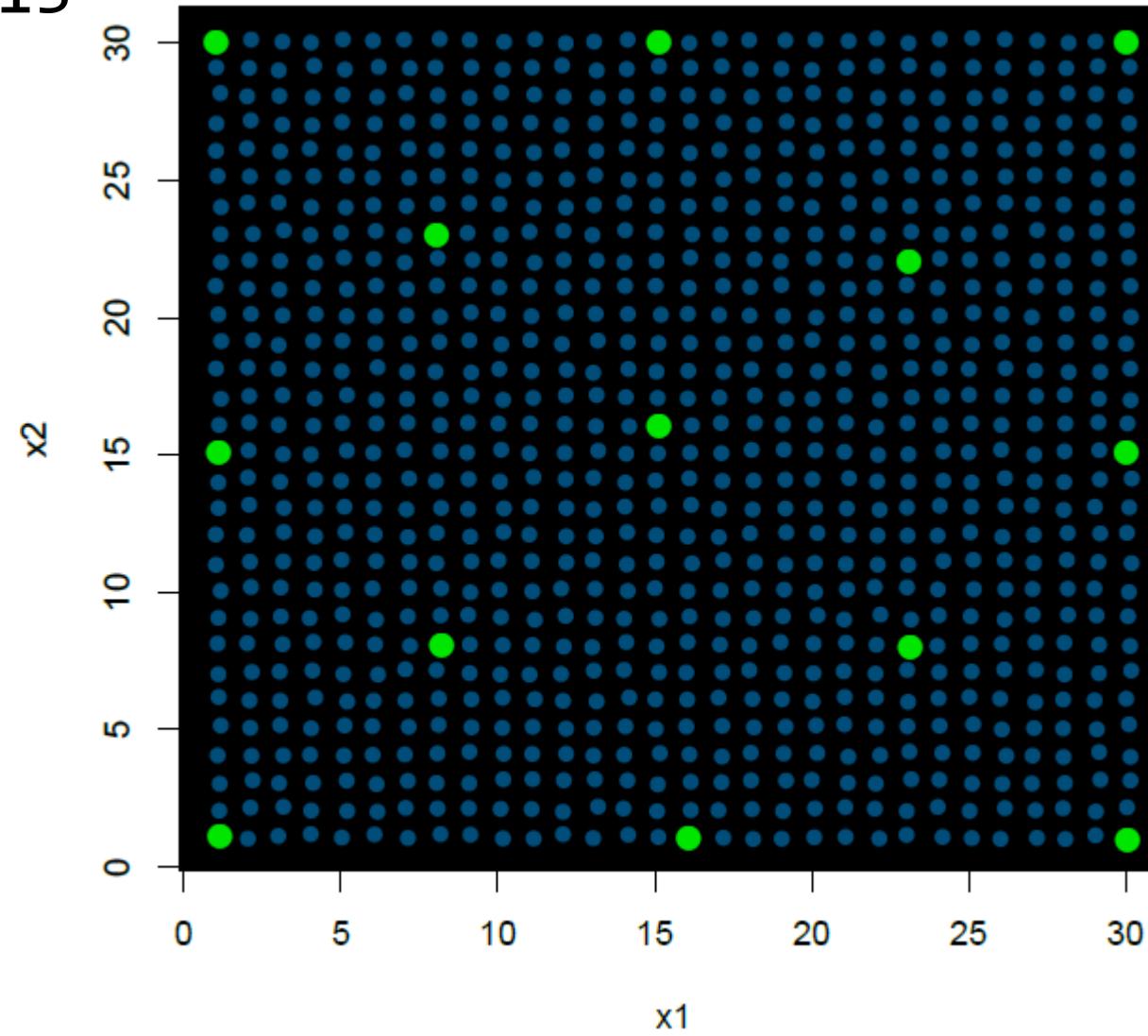


Calibration sampling

Sampling design

n = 13

Kennard-Stone sampling (KSS)

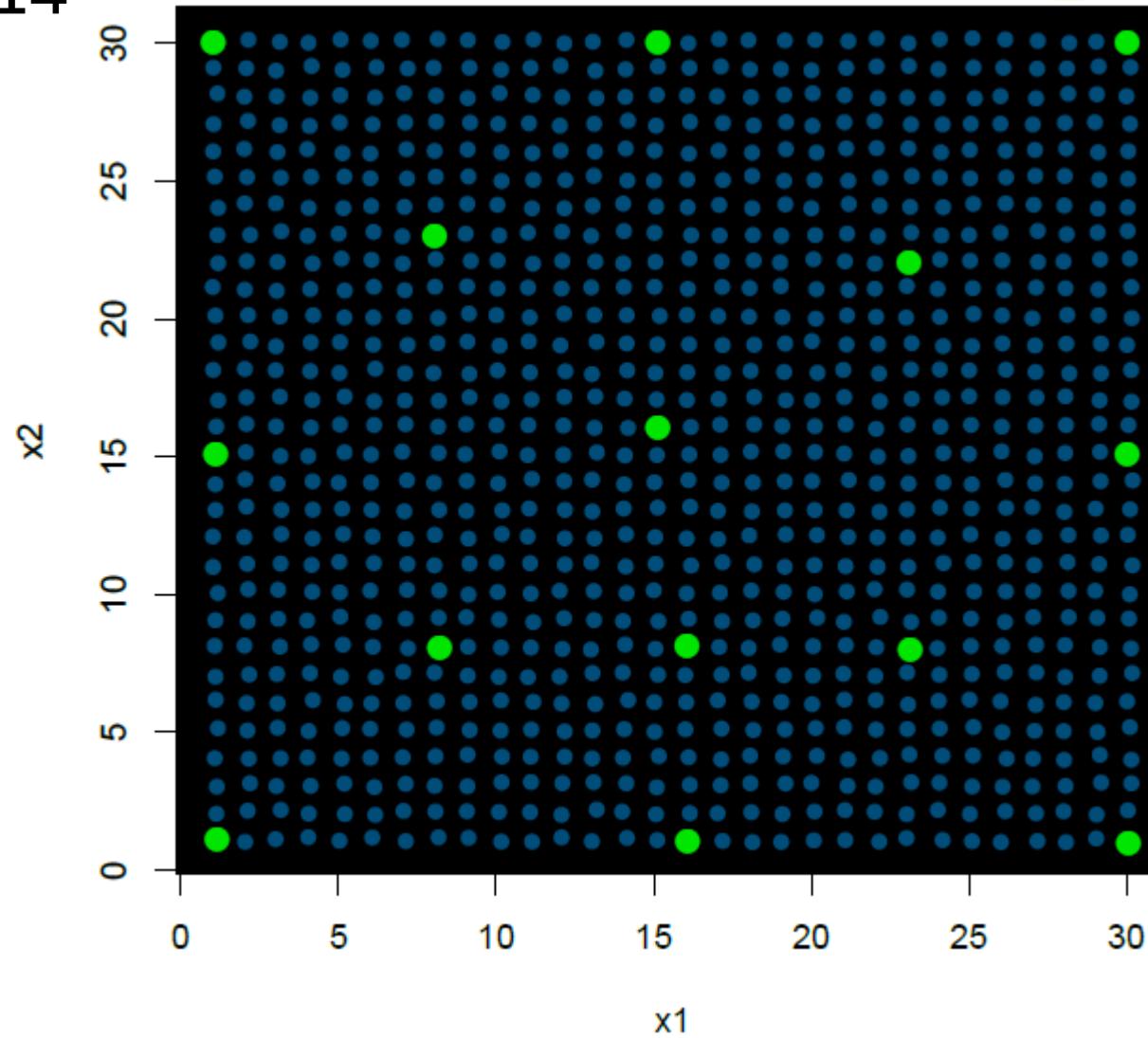


Calibration sampling

Sampling design

n = 14

Kennard-Stone sampling (KSS)

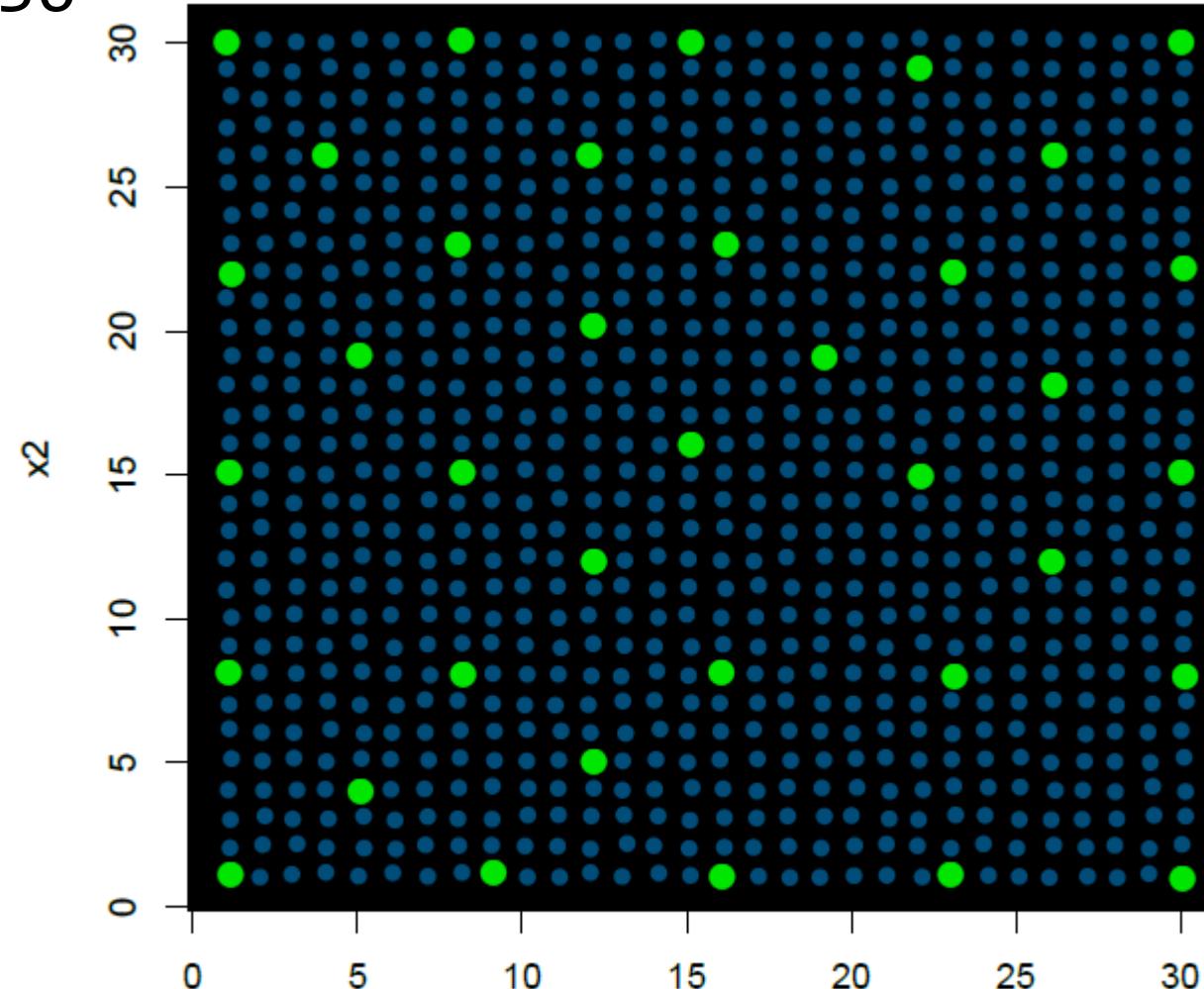


Calibration sampling

Sampling design

$n = 36$

Kennard-Stone sampling (KSS)



If there are outliers KSS will select some of them

K-means sampling (KMS)

The KMS algorithm (Naes, 1987) to select a (calibration) subset of m samples (S_{tr}) from a given set of N samples (S) works as follows:

1. Carry out a k -means clustering on the principal component scores and choose the number of resulting clusters to be equal to the number of desired calibration samples (m),
2. Select one sample from each cluster (usually the closest to the center of the cluster).

The **naes** function of the **prospectr** package can be used for KMS

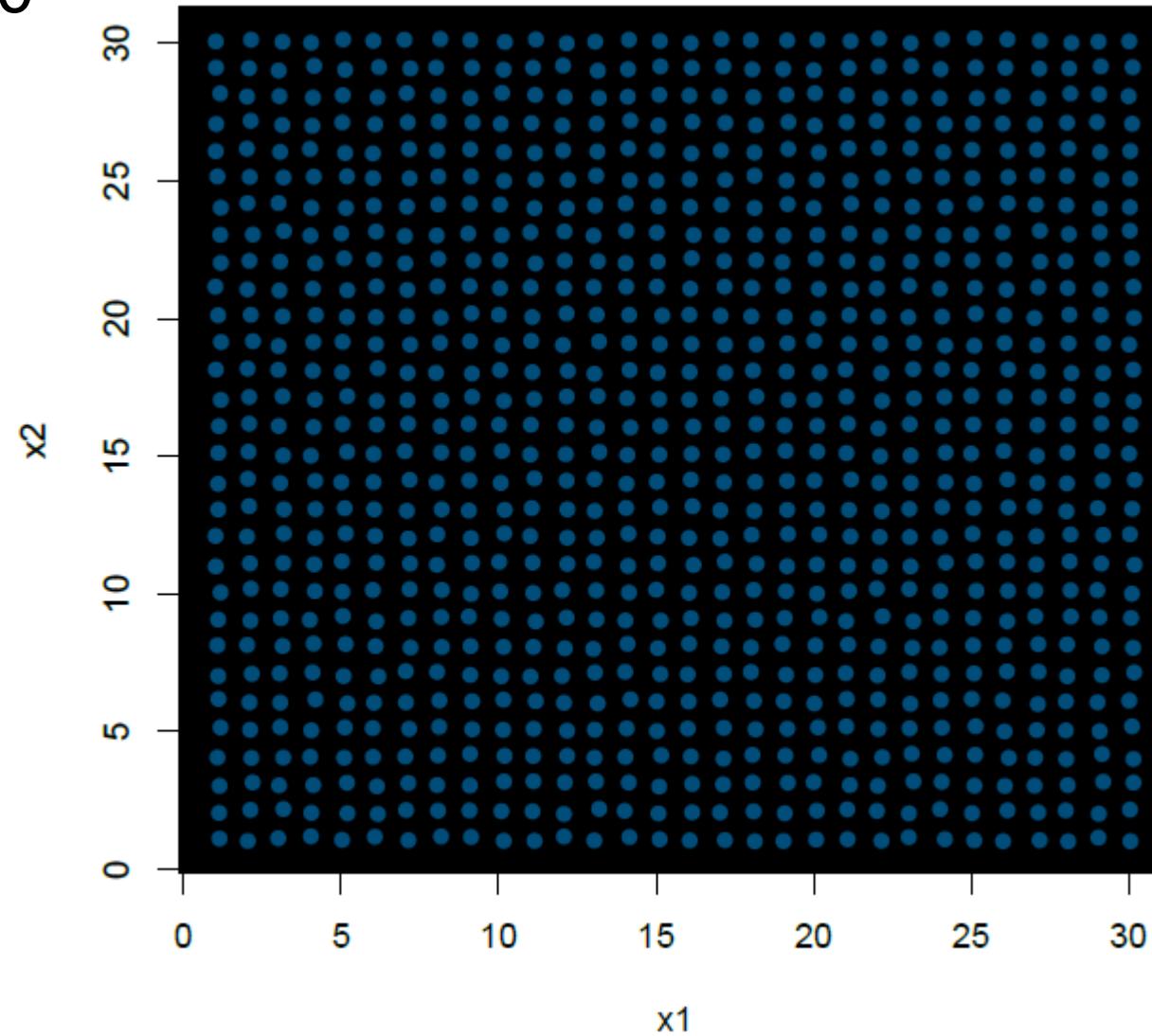
Naes, T., 1987. The design of calibration in near infra-red reflectance analysis by clustering. *Journal of Chemometrics* 1, 121-134.

Calibration sampling

Sampling design

$n = \emptyset$

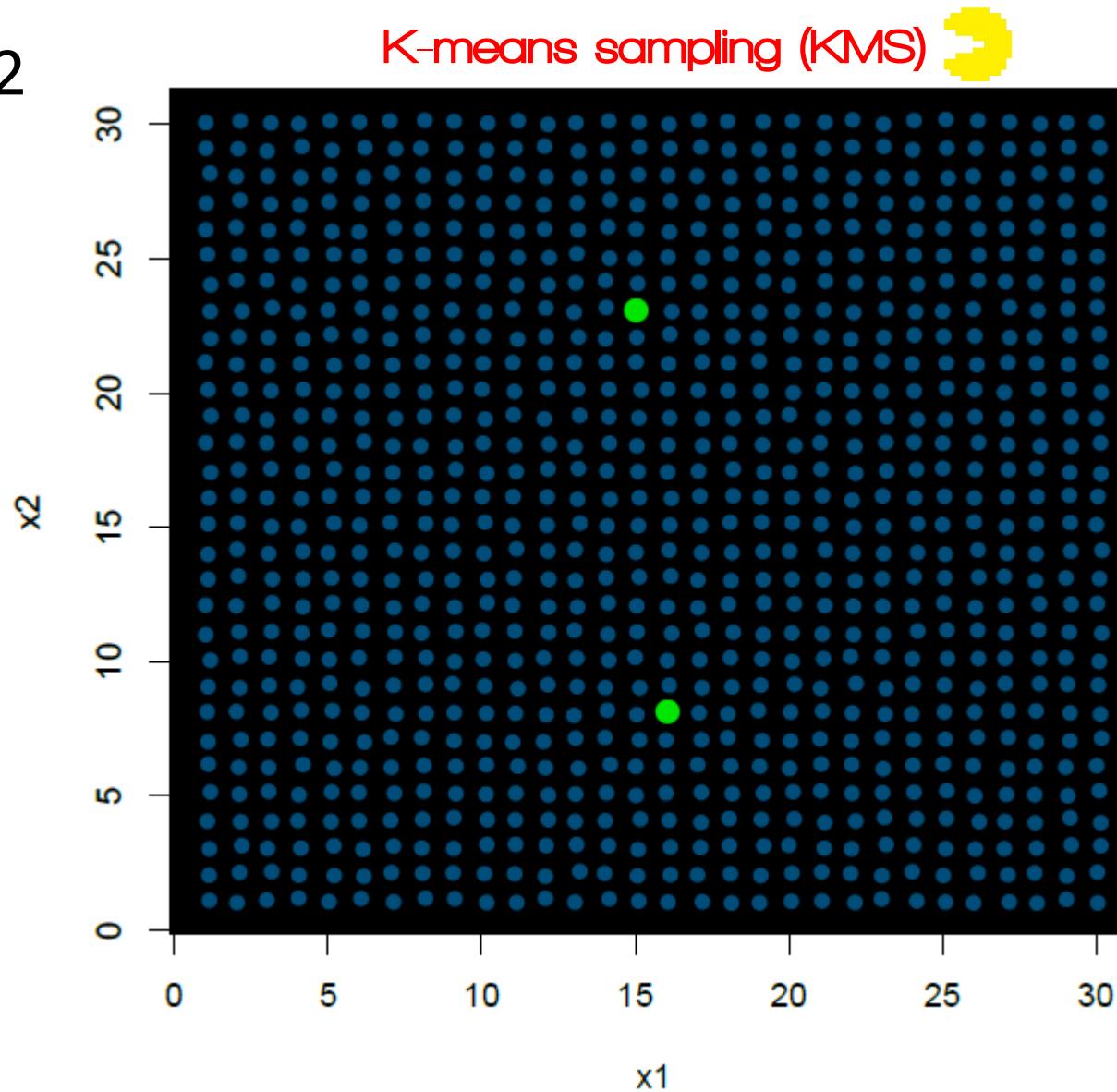
K-means sampling (KMS) 3



Calibration sampling

Sampling design

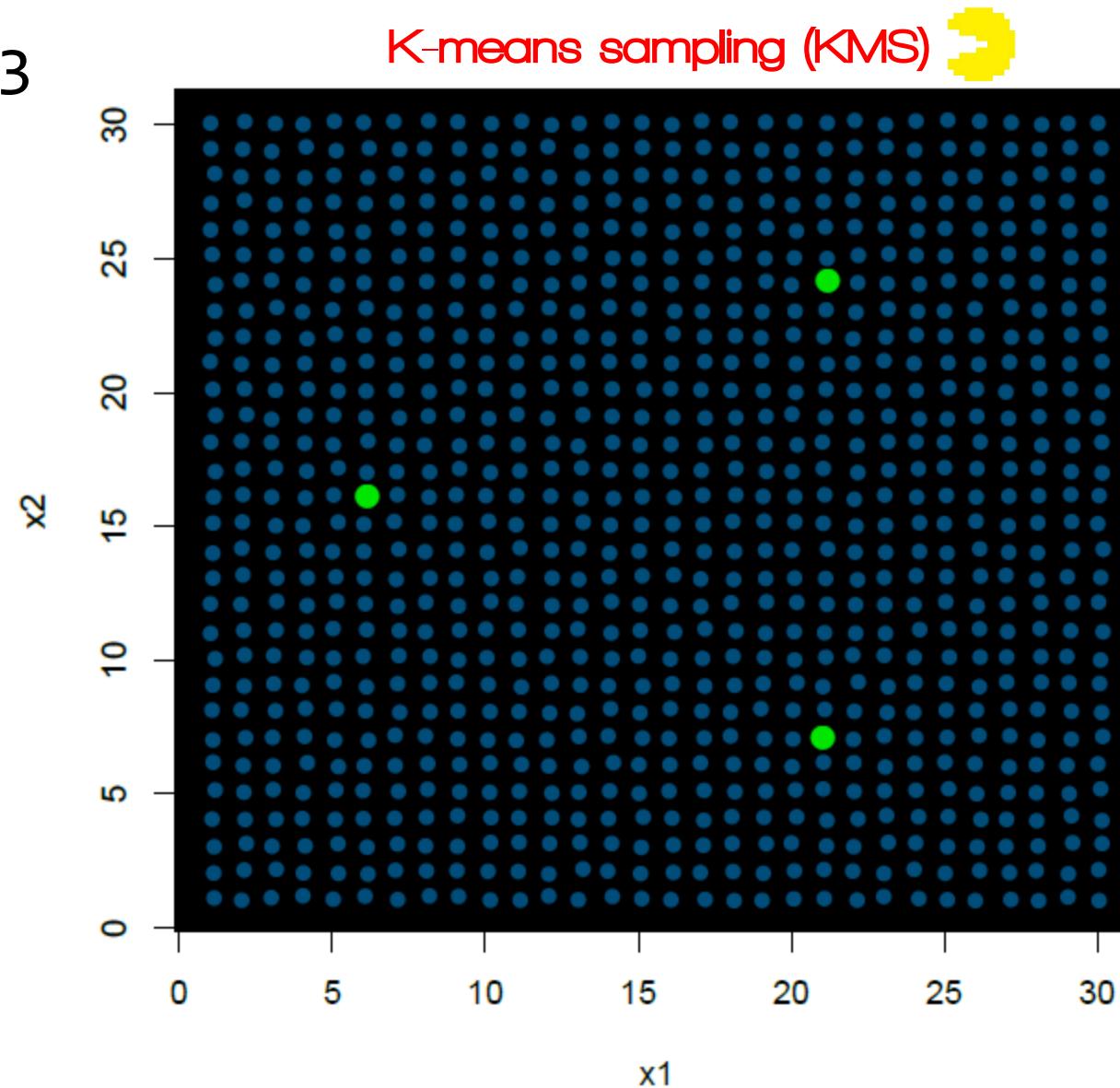
n = 2



Calibration sampling

Sampling design

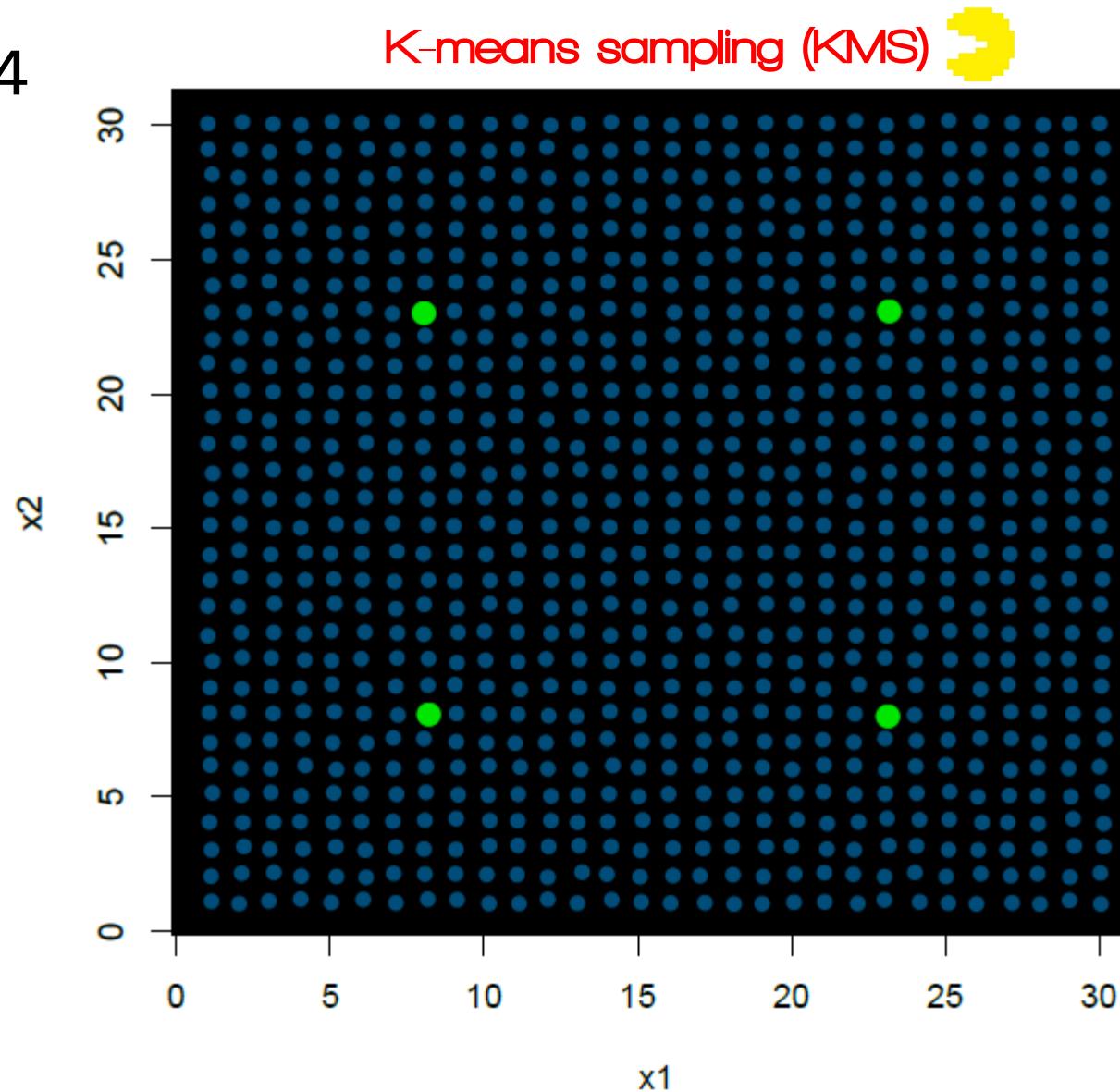
n = 3



Calibration sampling

Sampling design

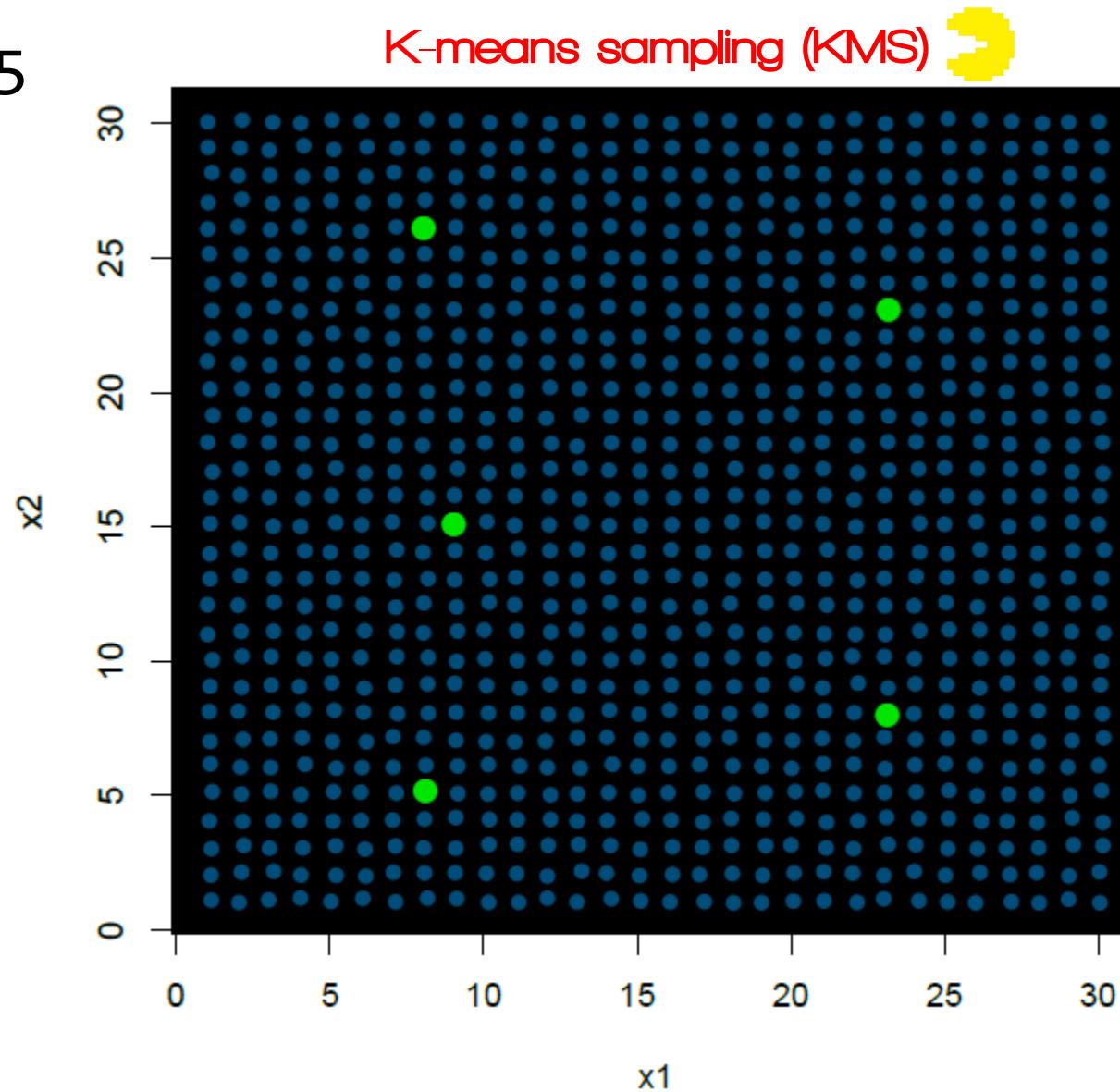
n = 4



Calibration sampling

Sampling design

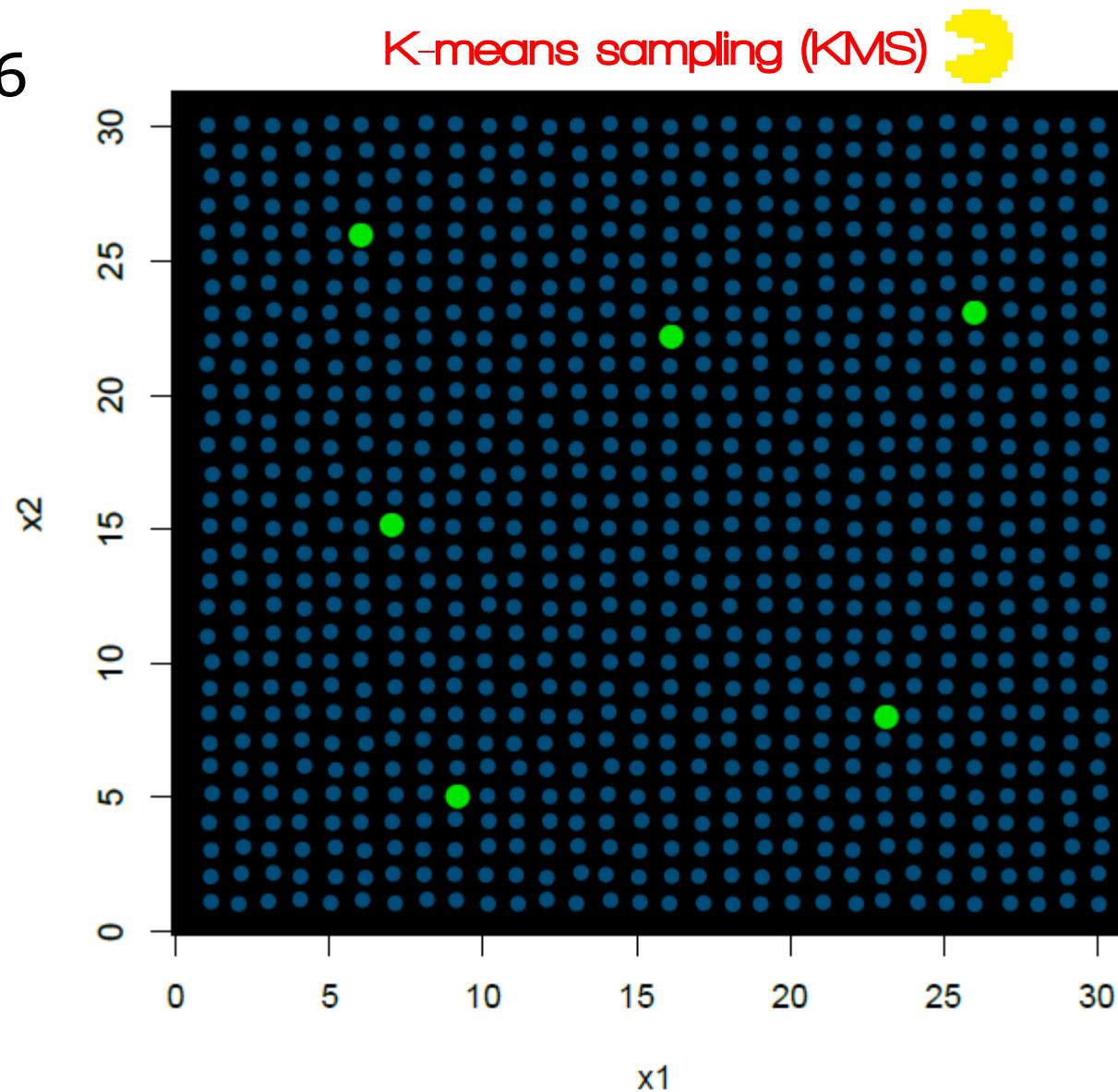
n = 5



Calibration sampling

Sampling design

n = 6

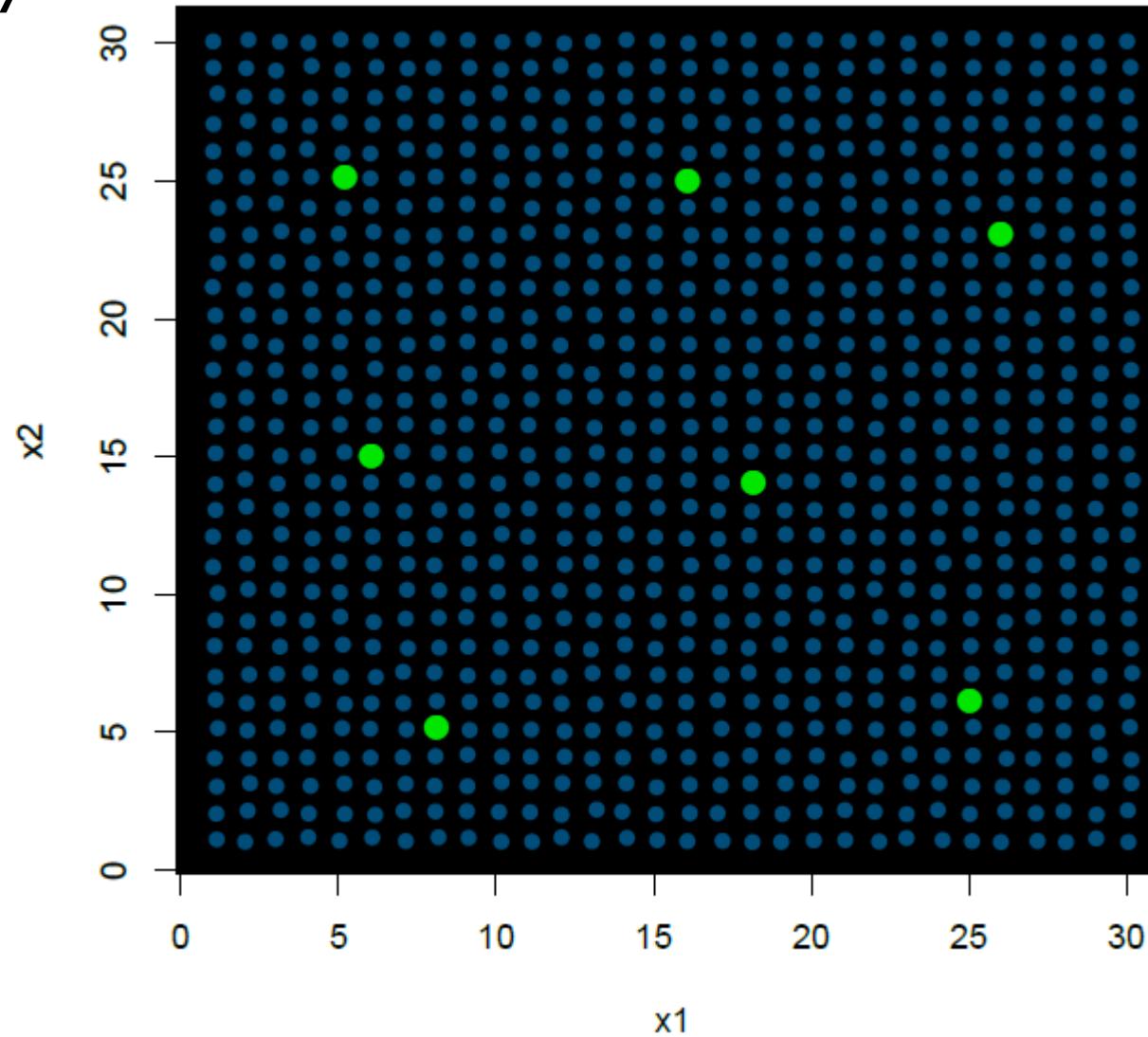


Calibration sampling

Sampling design

n = 7

K-means sampling (KMS) 3

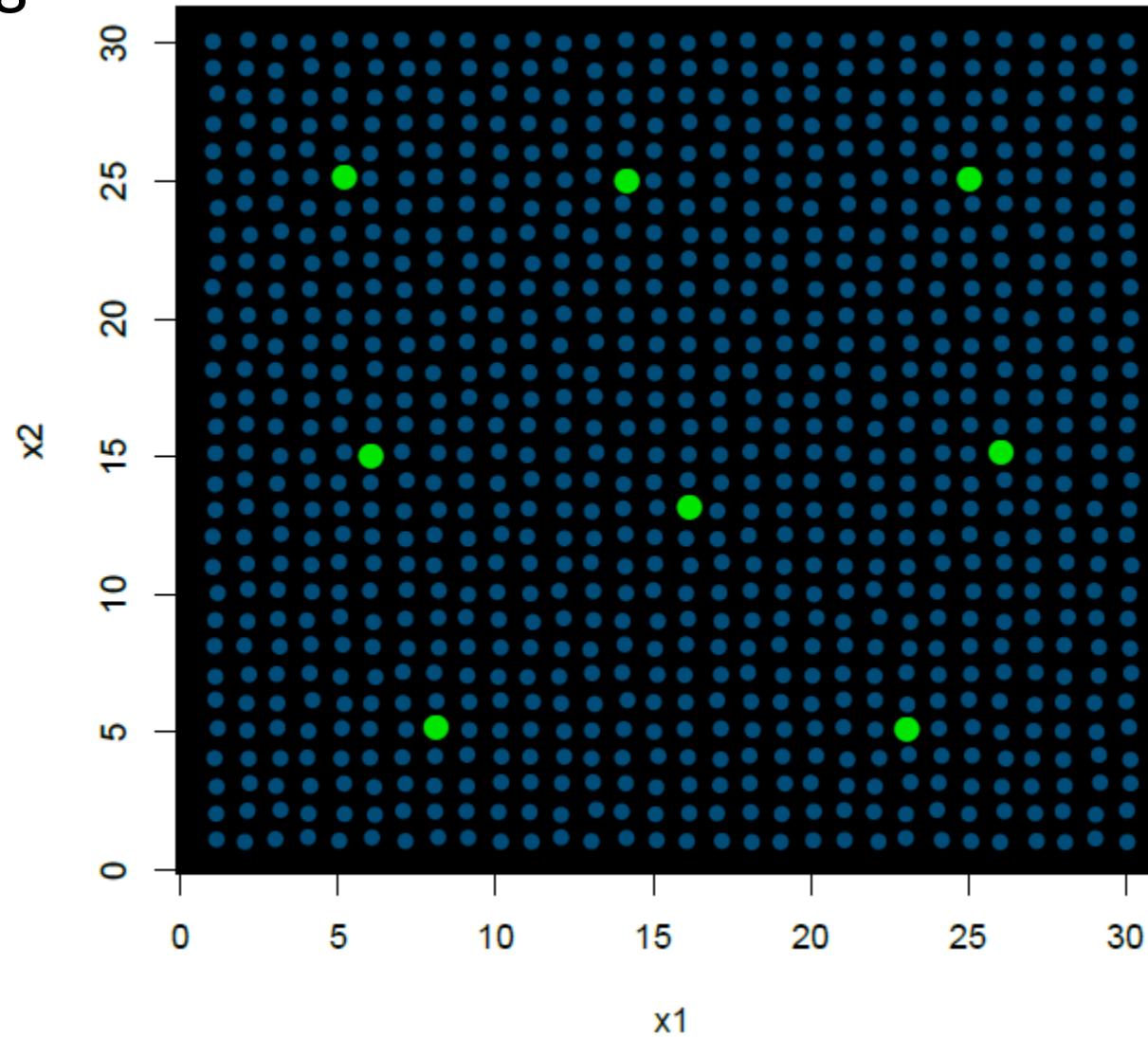


Calibration sampling

Sampling design

n = 8

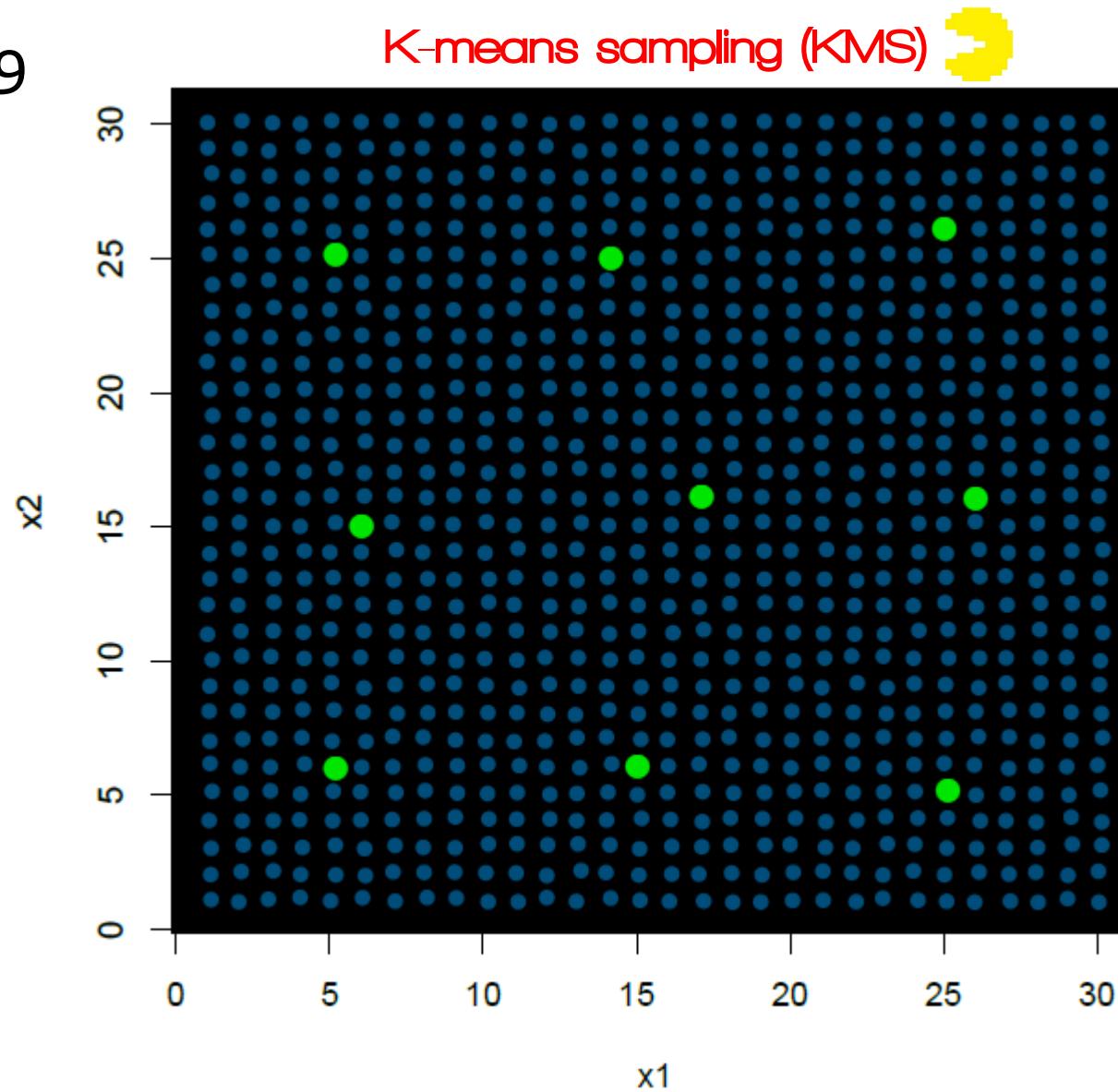
K-means sampling (KMS) 3



Calibration sampling

Sampling design

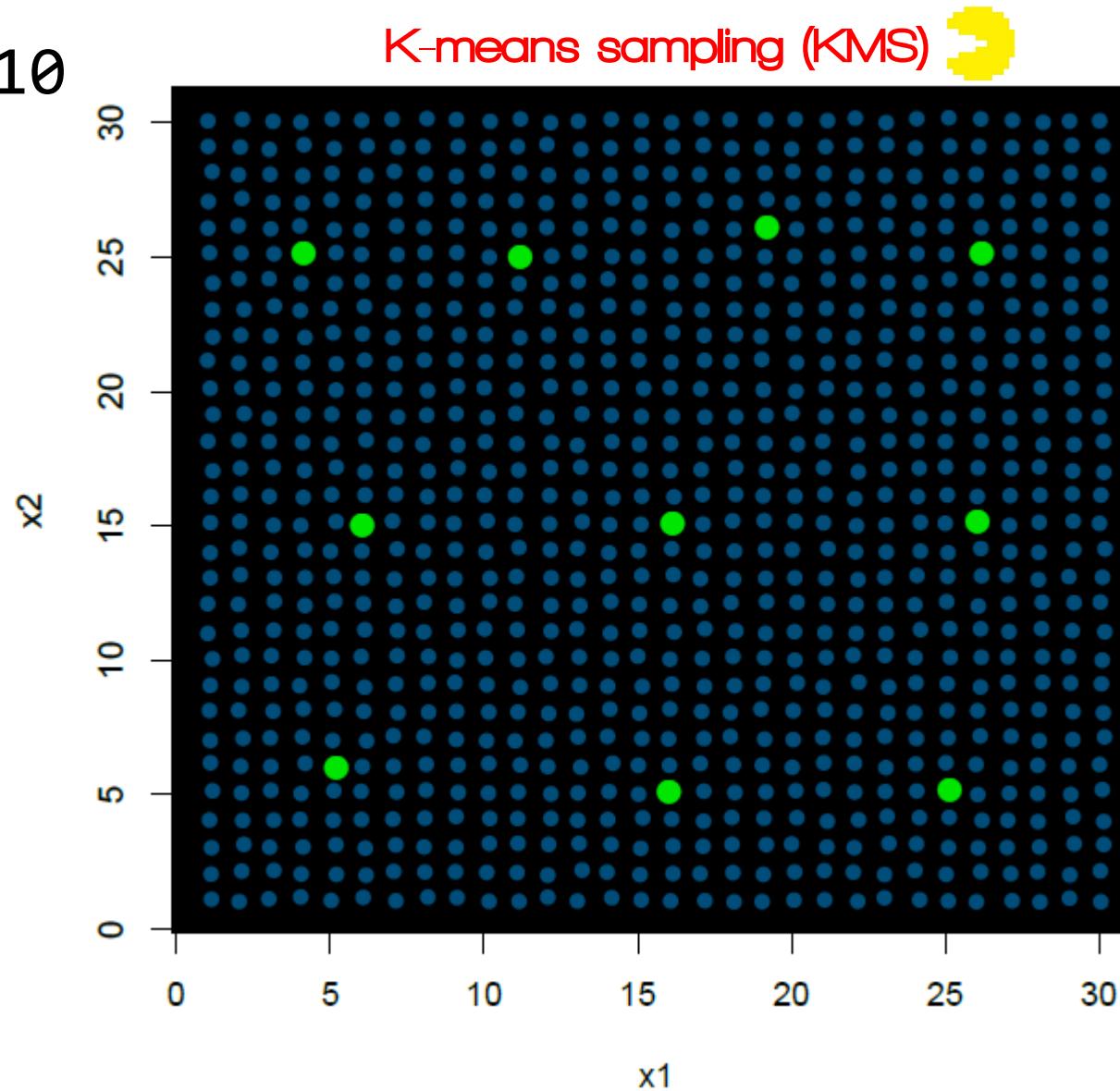
n = 9



Calibration sampling

Sampling design

n = 10

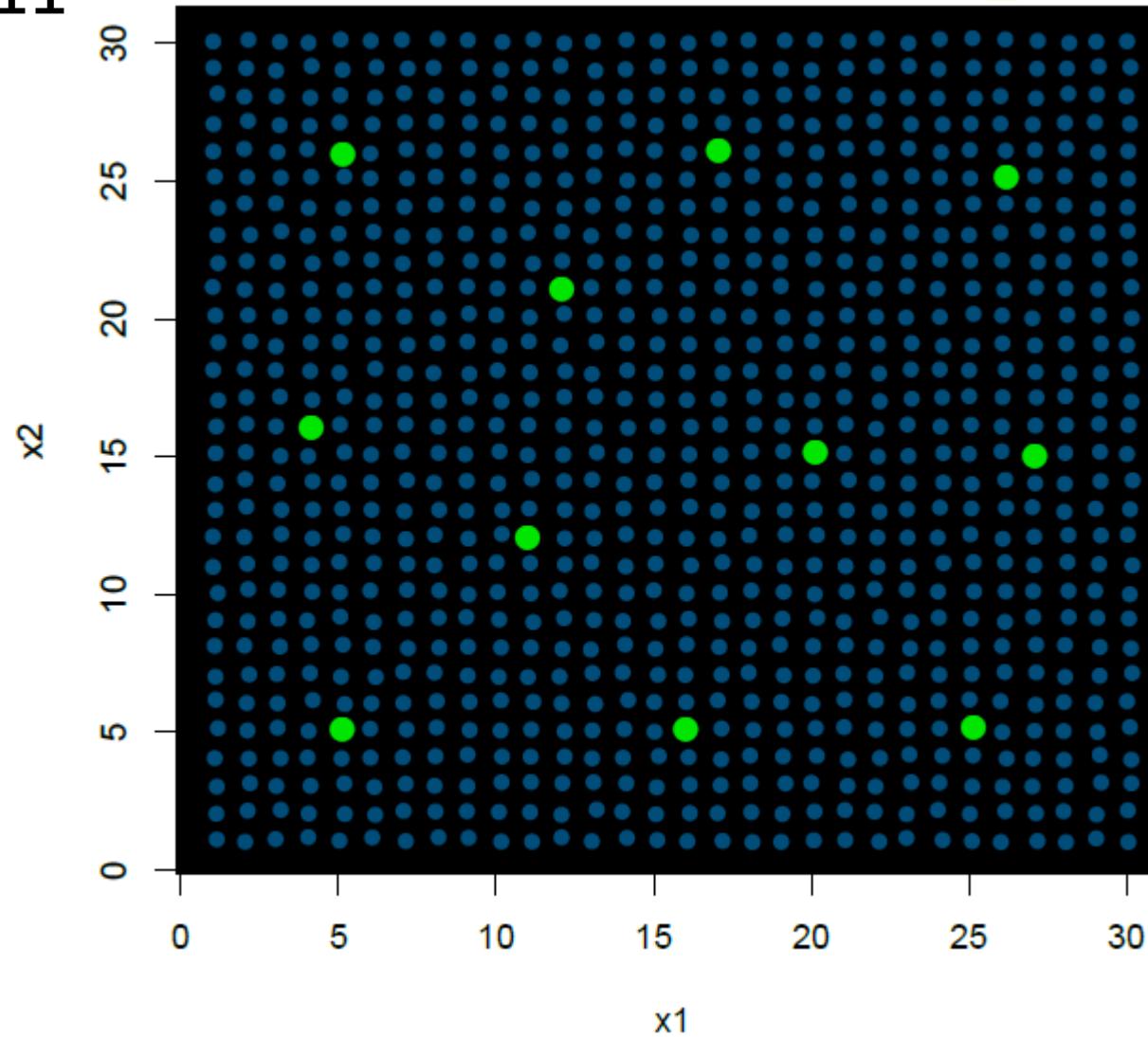


Calibration sampling

Sampling design

n = 11

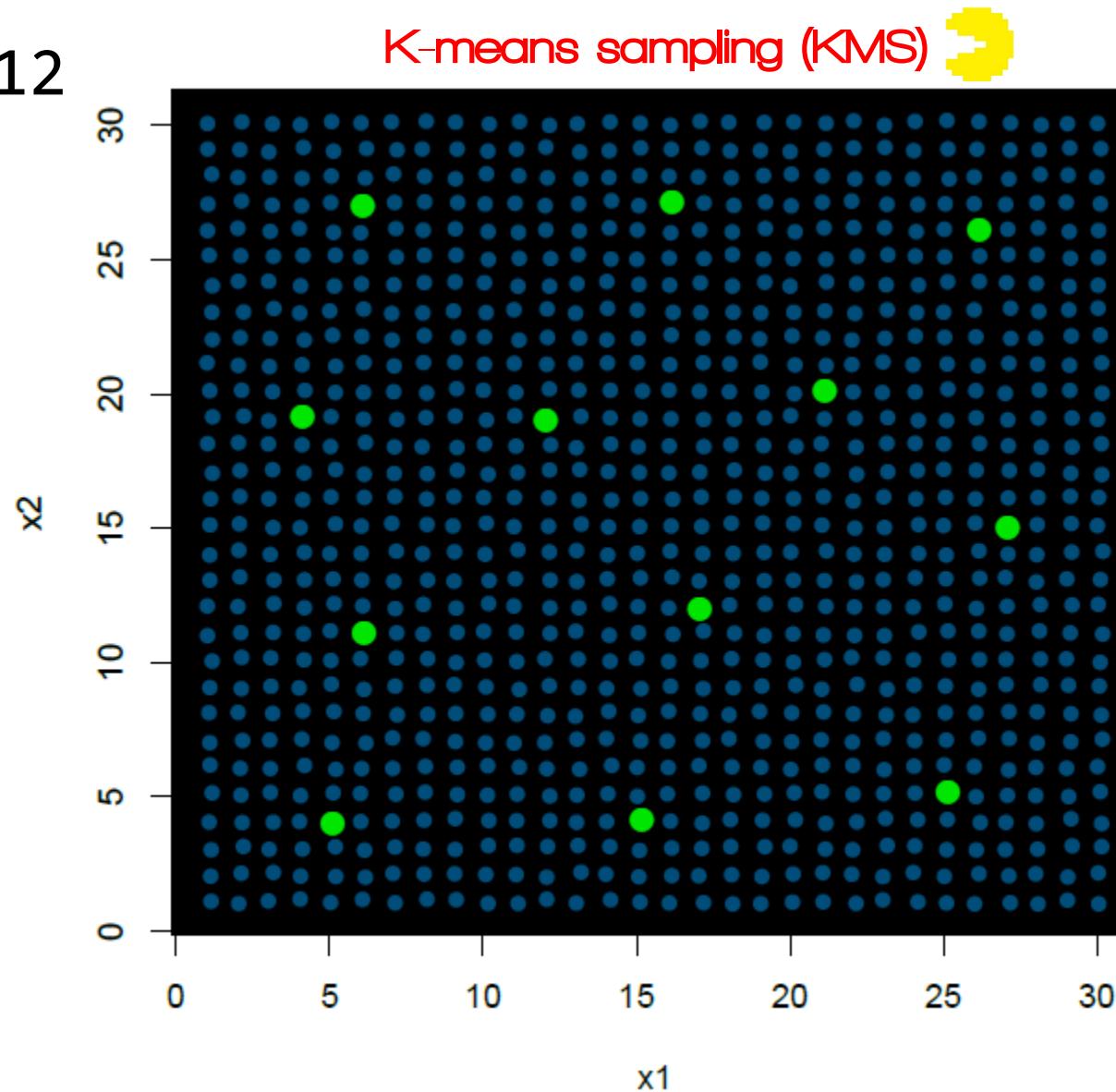
K-means sampling (KMS) 3



Calibration sampling

Sampling design

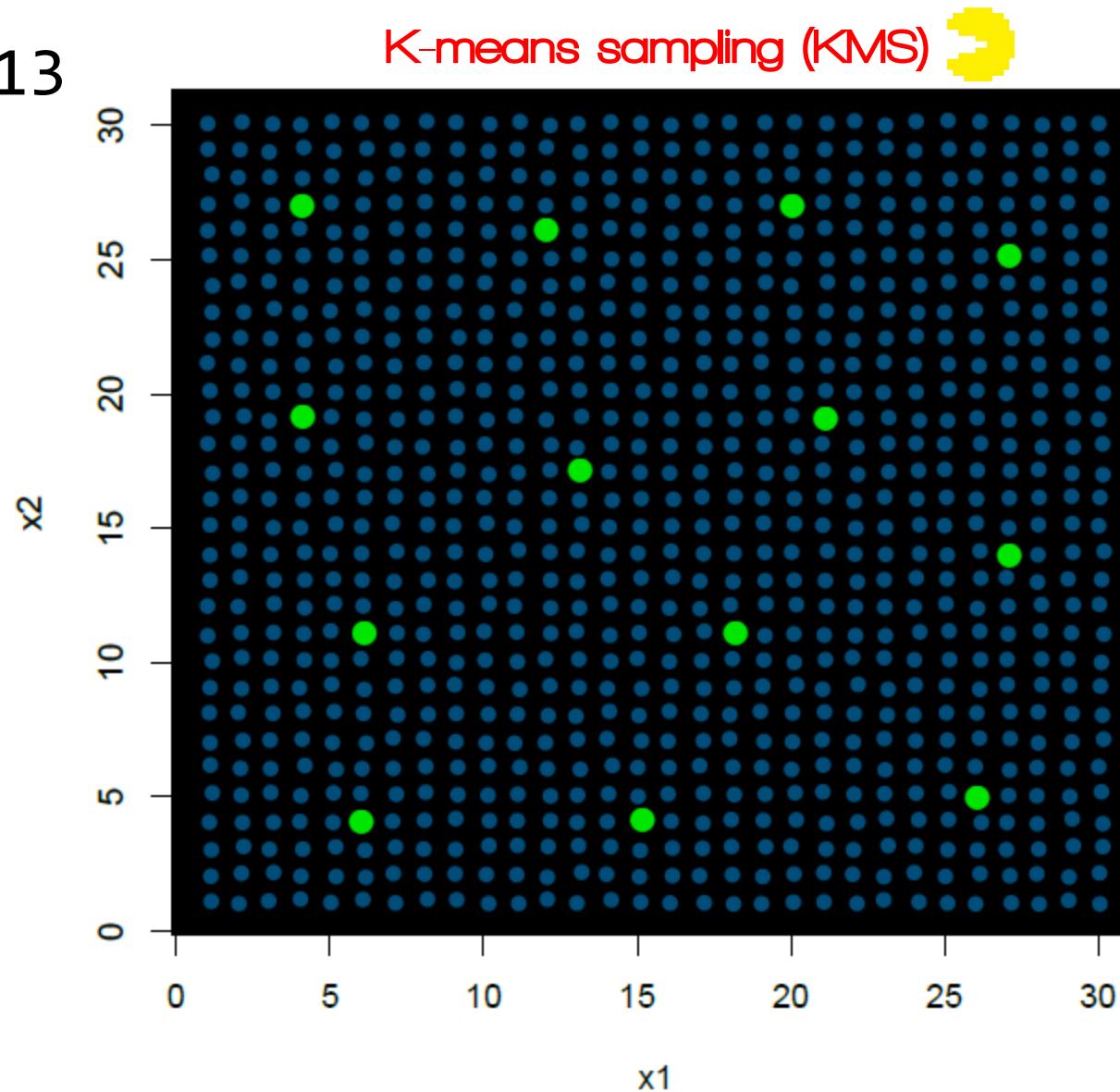
n = 12



Calibration sampling

Sampling design

n = 13

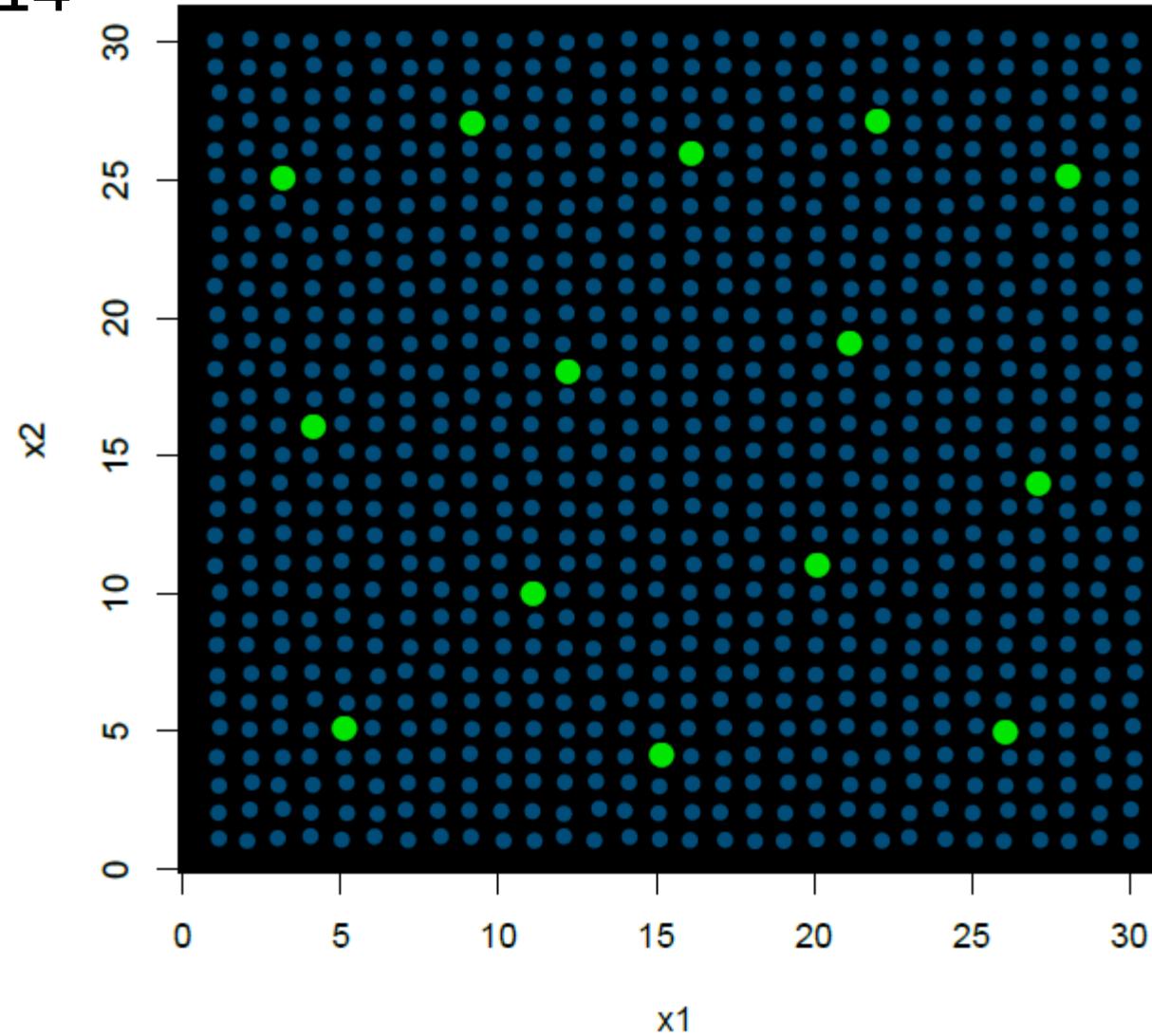


Calibration sampling

Sampling design

n = 14

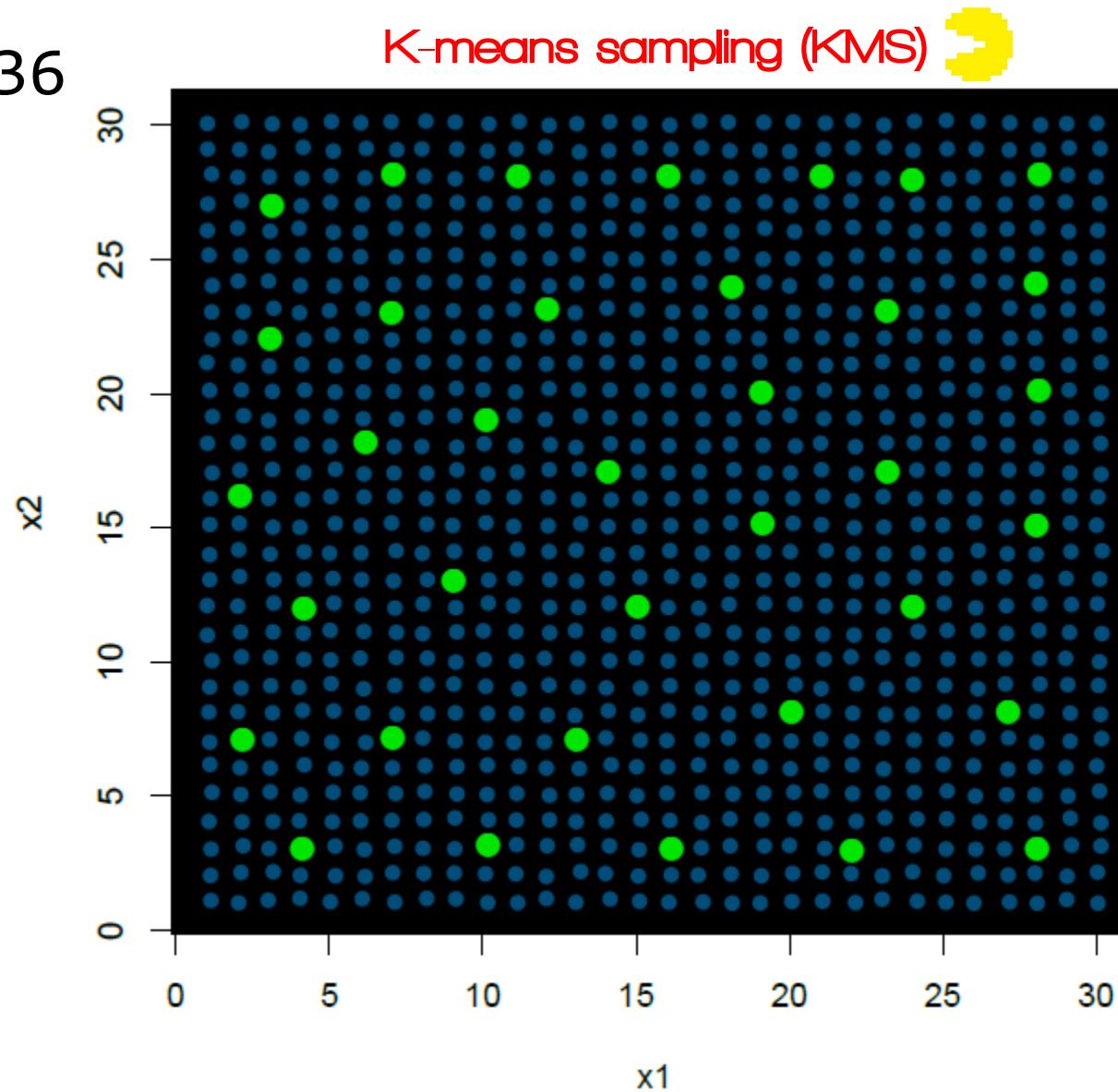
K-means sampling (KMS) 3



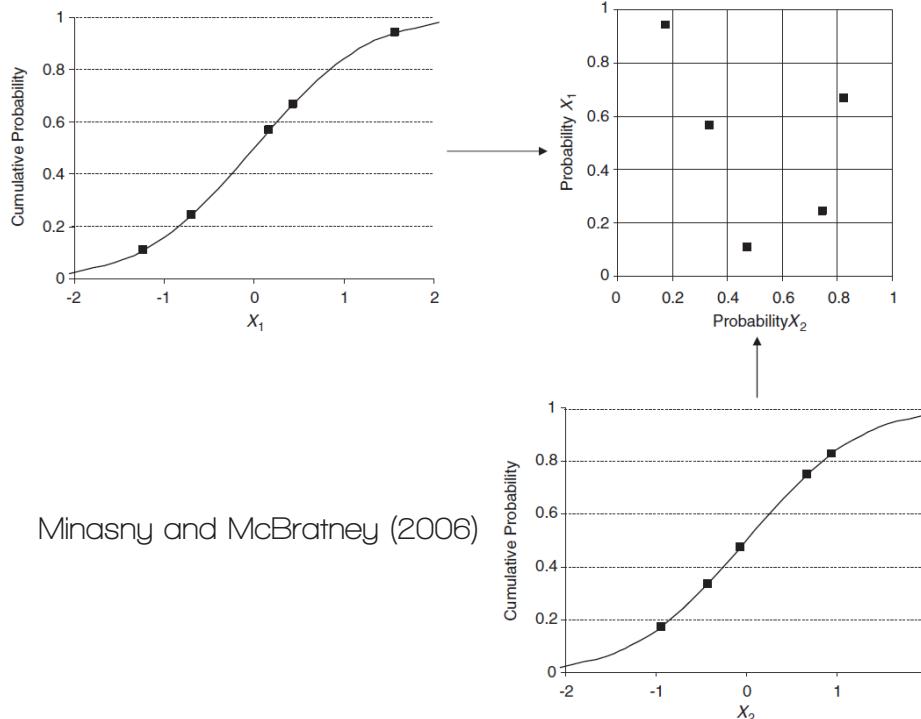
Calibration sampling

Sampling design

n = 36



Conditioned Latin hypercube sampling (cLHS)



The cLHS attempts to cover the multidimensional distribution corresponding to a set of variables by using a stratified random sampling.

The **clhs** function of the **clhs** package can be used for cLHS

cLHS does not maximize the dissimilarity between samples

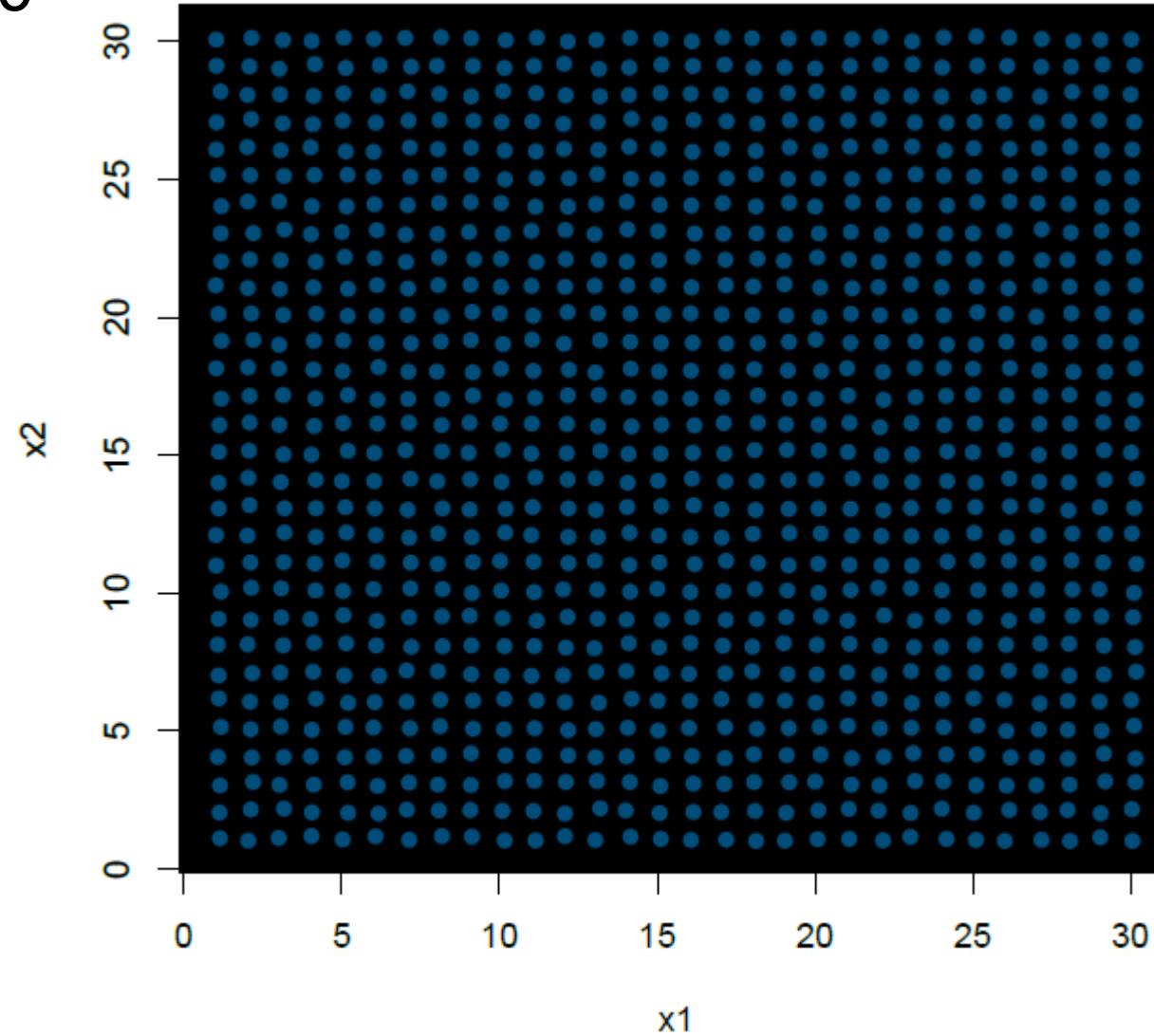
Minasny, B., & McBratney, A. B. (2006). A conditioned Latin hypercube method for sampling in the presence of ancillary information. *Computers & Geosciences*, 32(9), 1378-1388.

Calibration sampling

Sampling design

$n = \emptyset$

contioned Latin Hypercube sampling (cLHS)

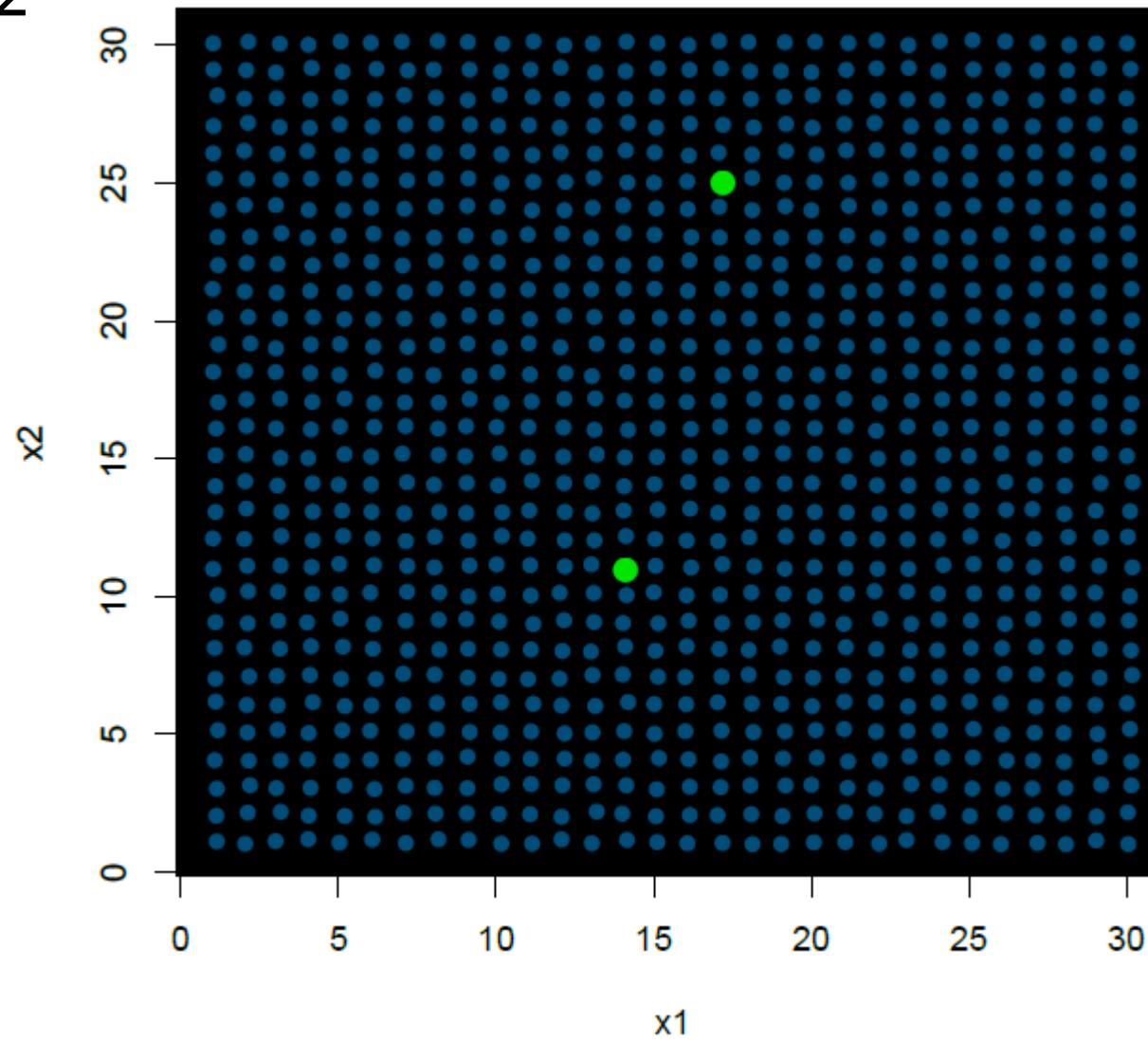


Calibration sampling

Sampling design

$n = 2$

contioned Latin Hypercube sampling (cLHS)

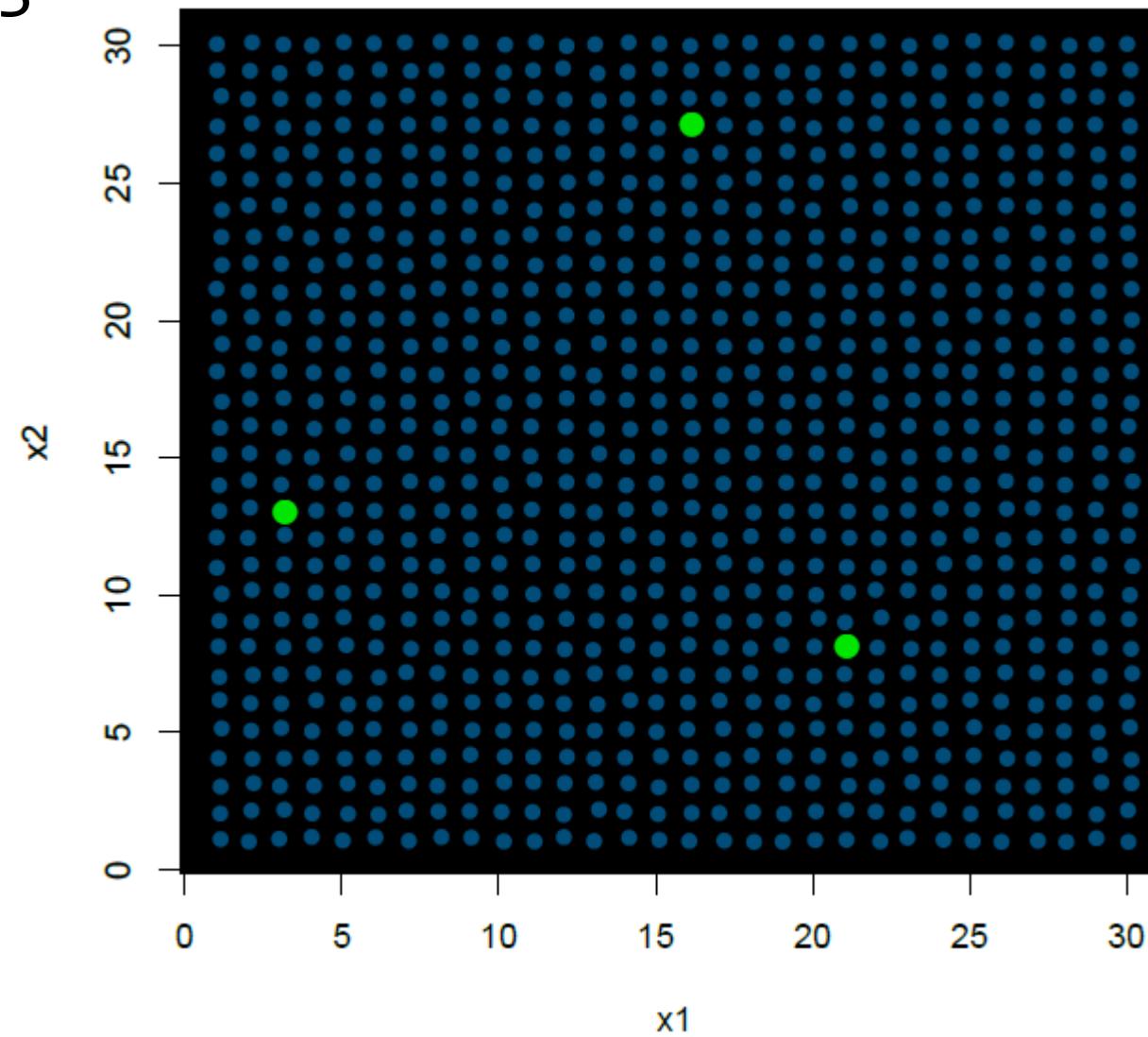


Calibration sampling

Sampling design

n = 3

contioned Latin Hypercube sampling (cLHS)

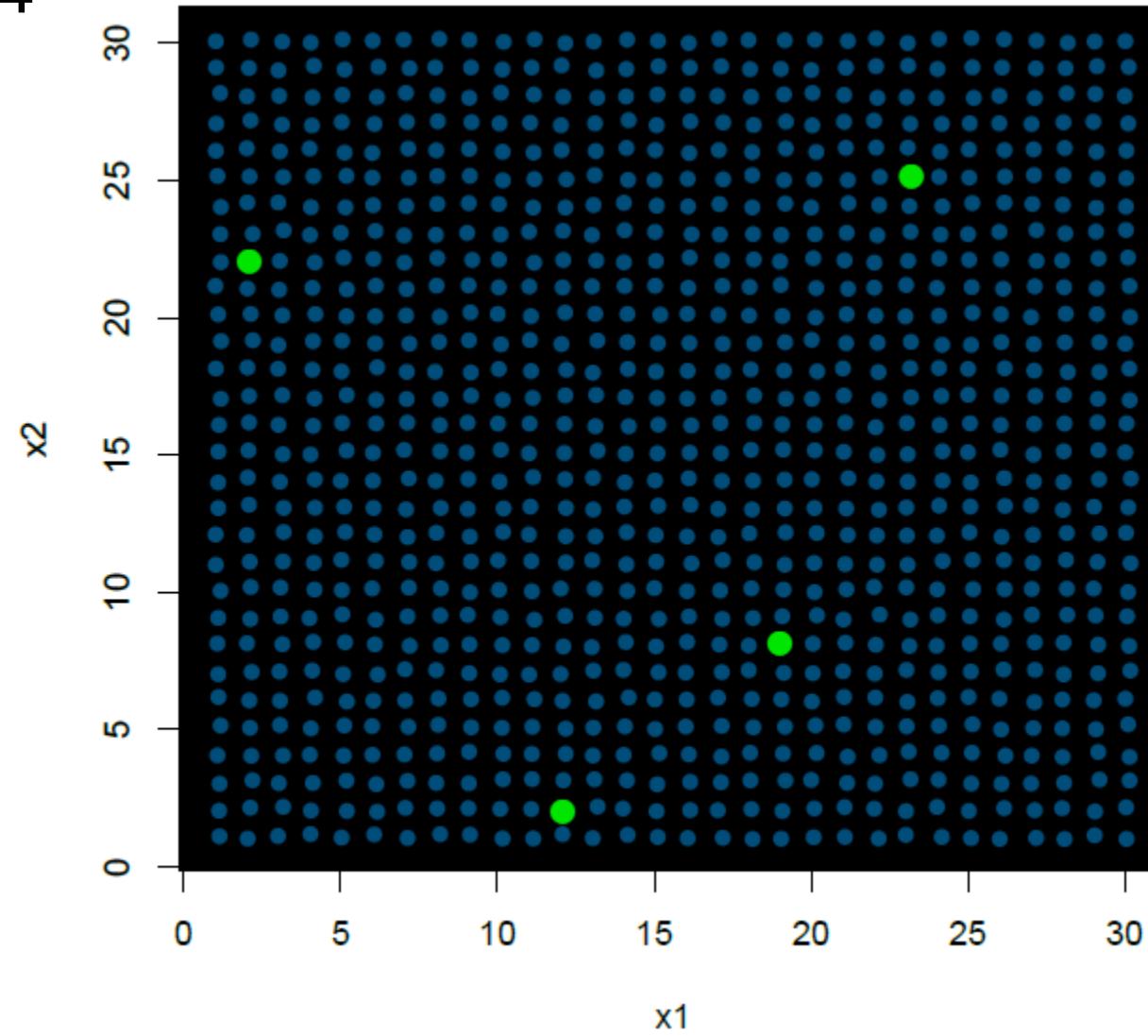


Calibration sampling

Sampling design

n = 4

contioned Latin Hypercube sampling (cLHS)

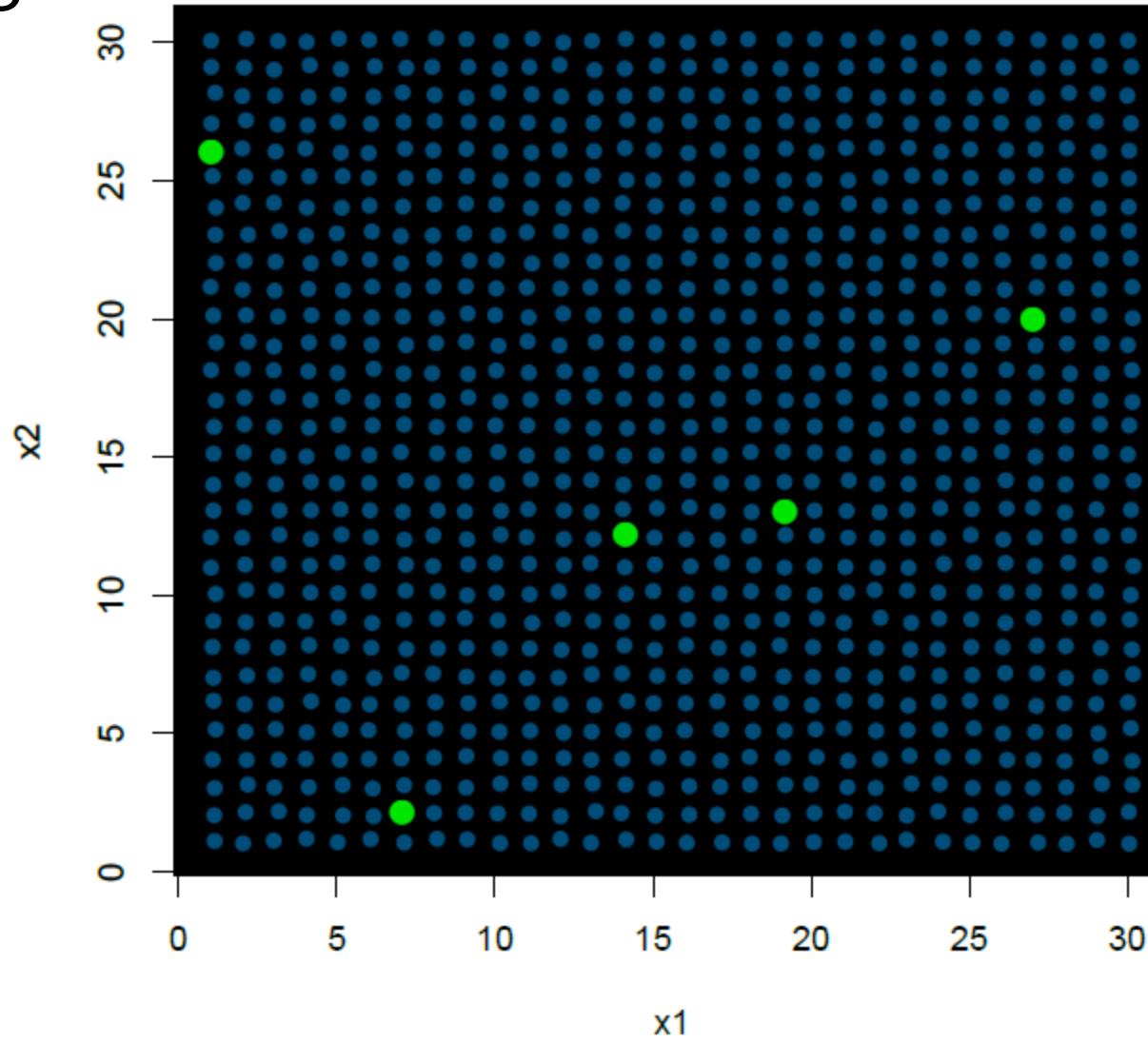


Calibration sampling

Sampling design

n = 5

contioned Latin Hypercube sampling (cLHS)

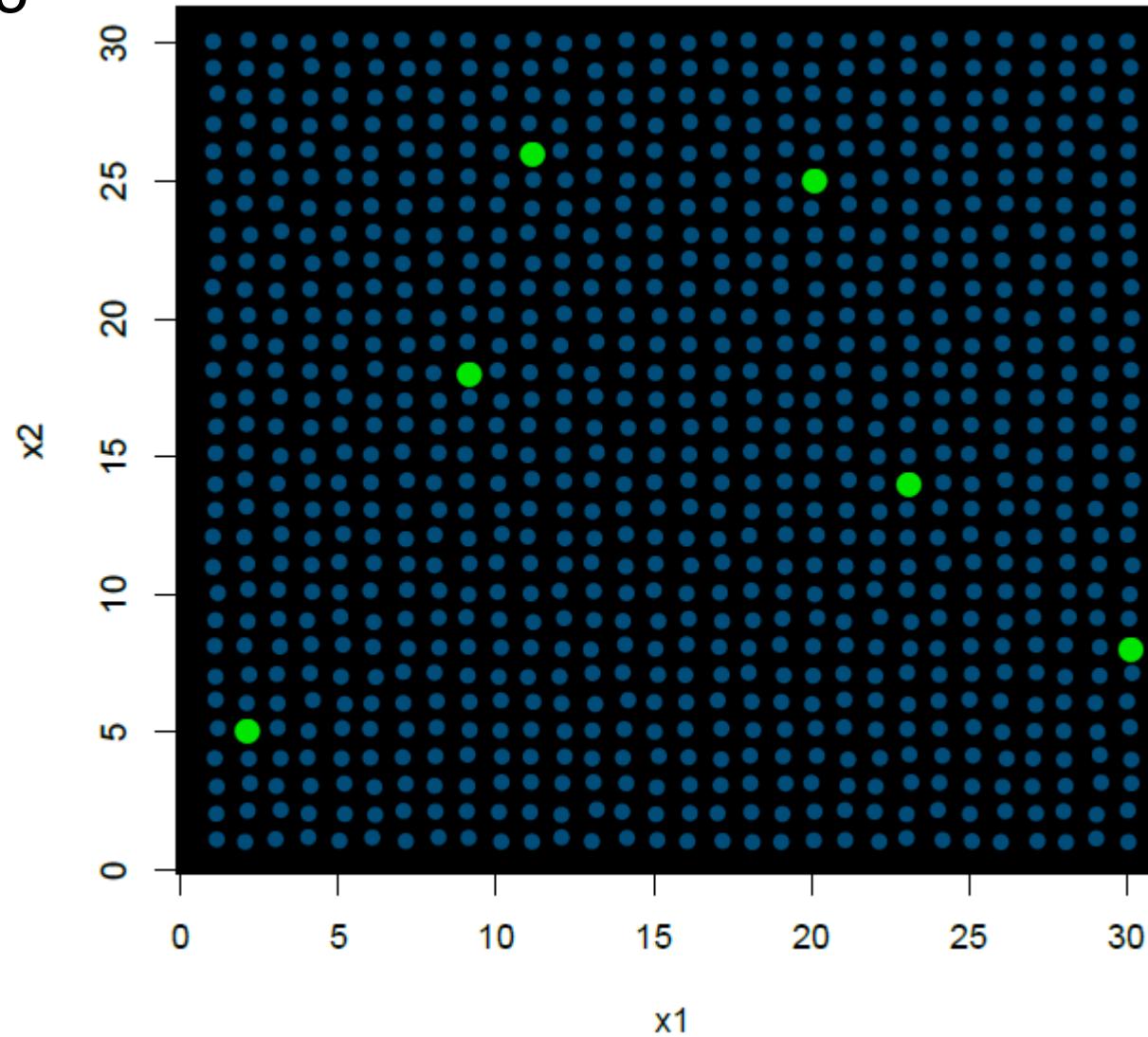


Calibration sampling

Sampling design

n = 6

contioned Latin Hypercube sampling (cLHS)

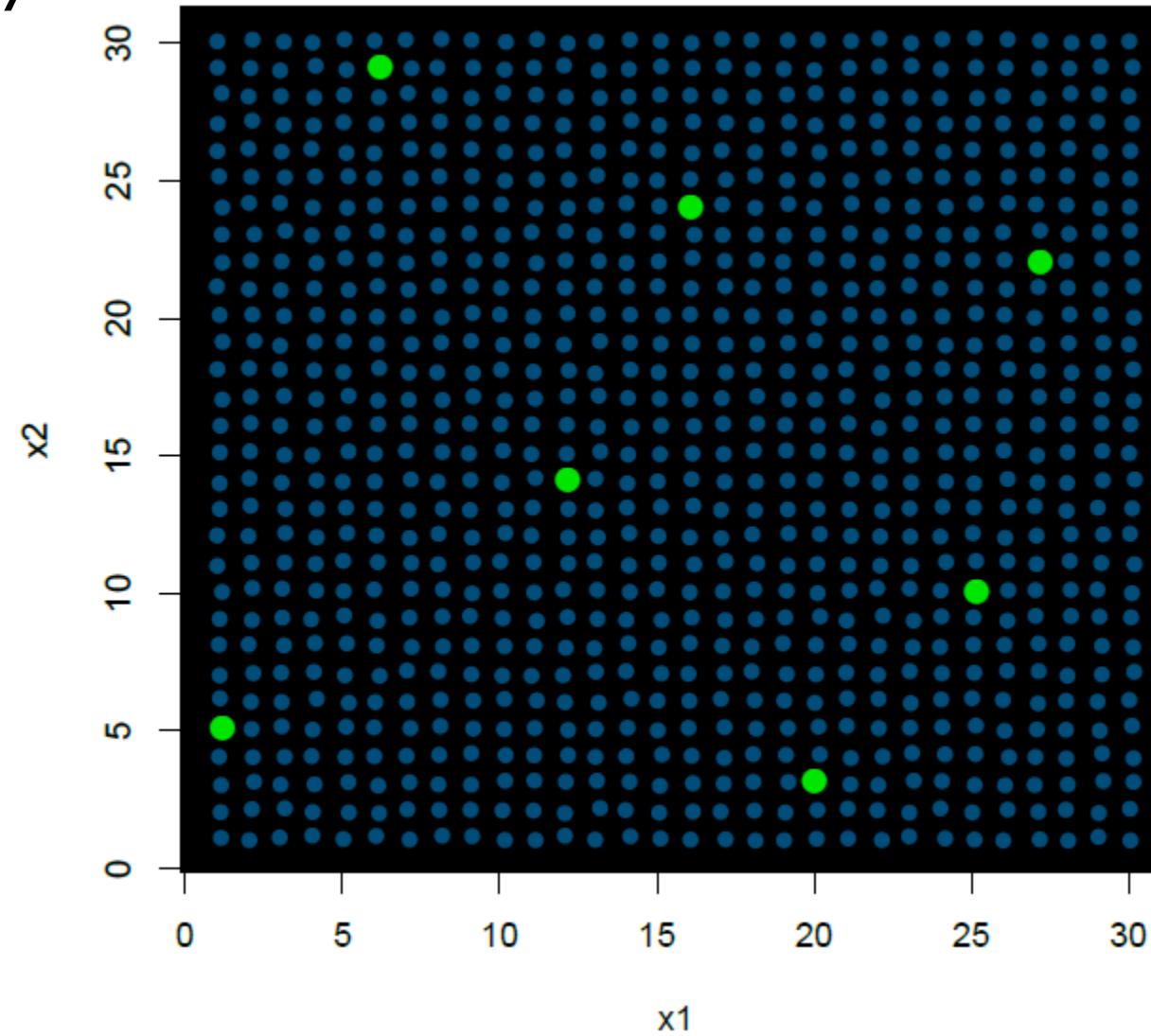


Calibration sampling

Sampling design

n = 7

contioned Latin Hypercube sampling (cLHS)

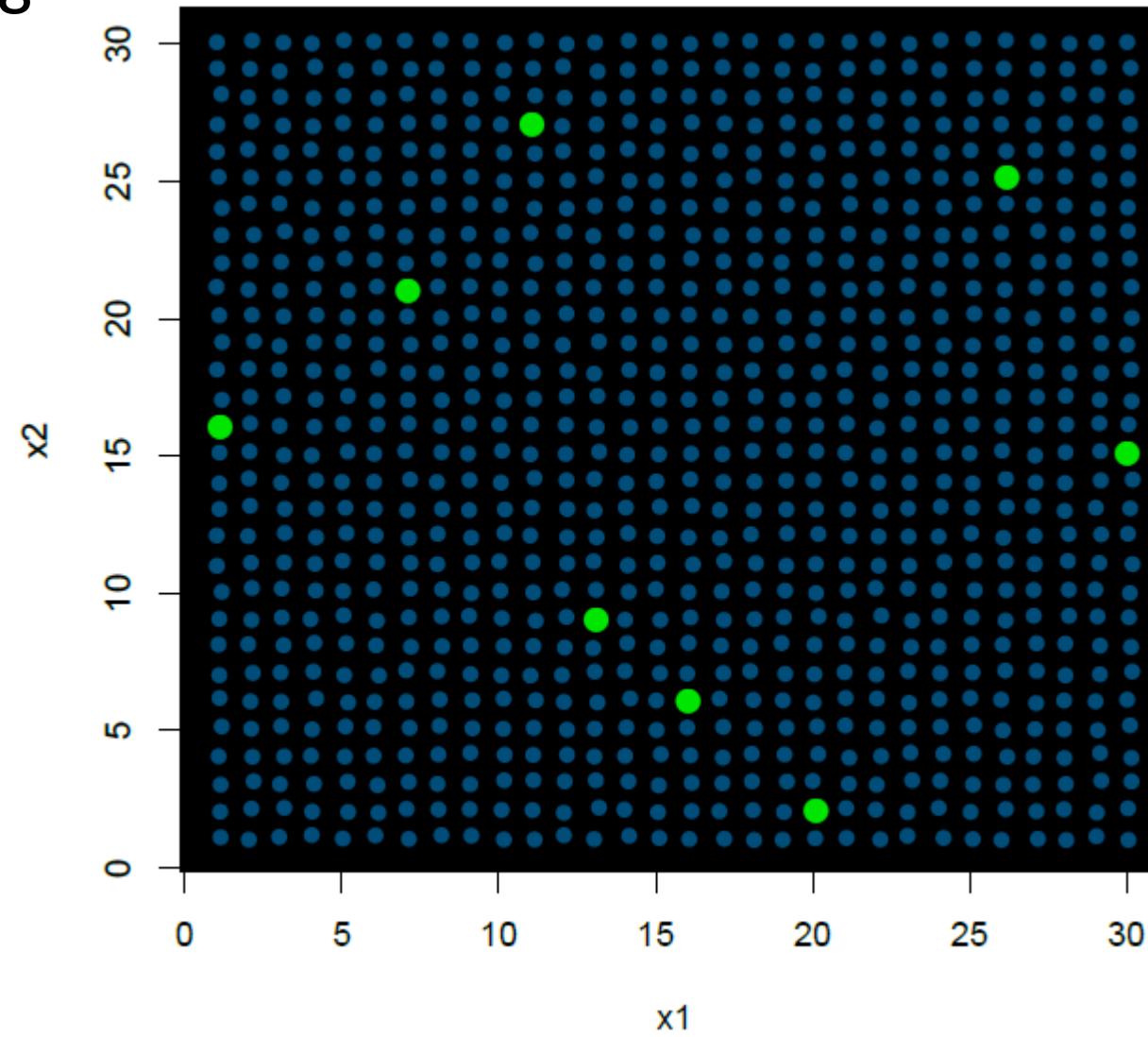


Calibration sampling

Sampling design

n = 8

contioned Latin Hypercube sampling (cLHS)

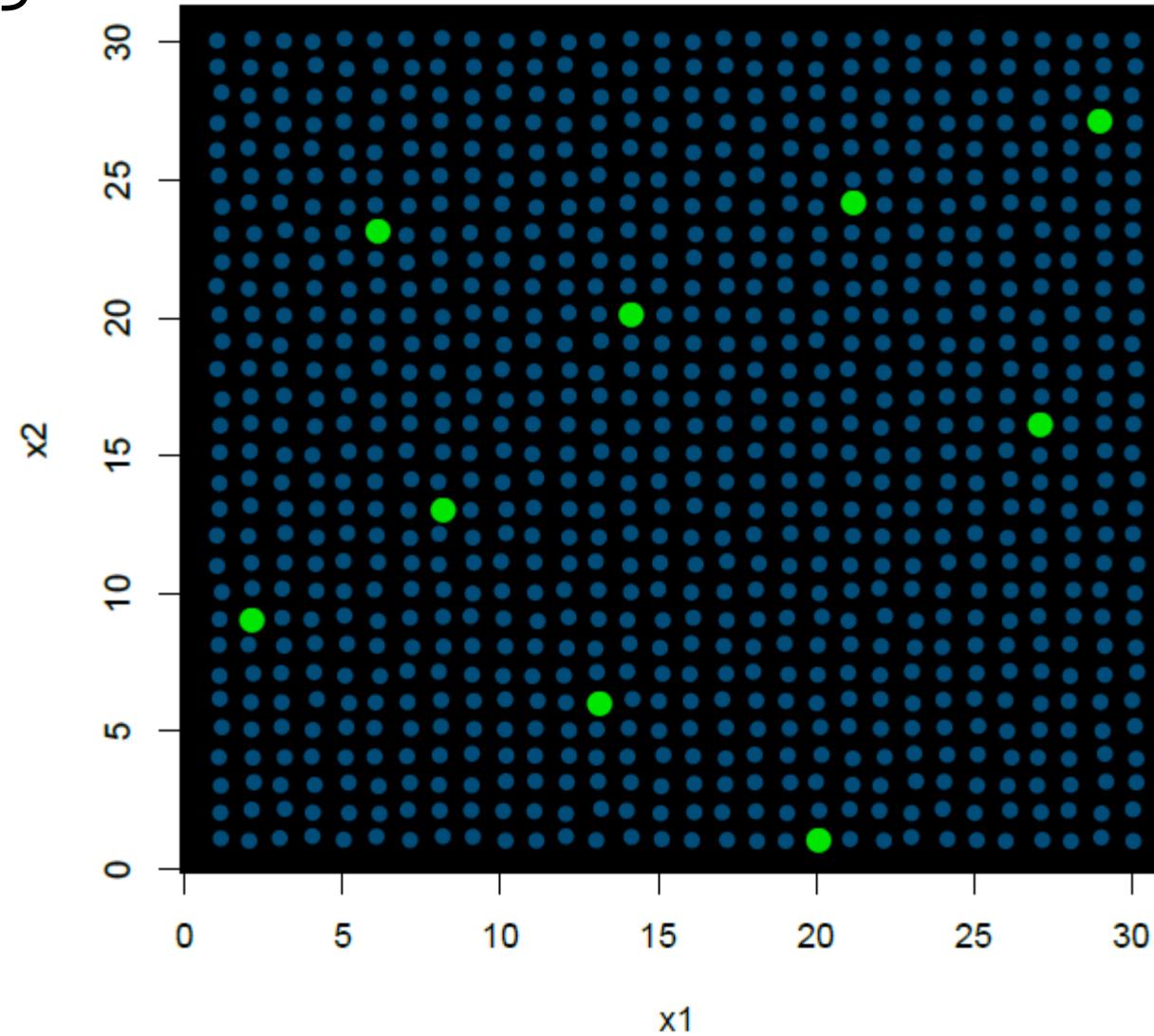


Calibration sampling

Sampling design

n = 9

contioned Latin Hypercube sampling (cLHS)

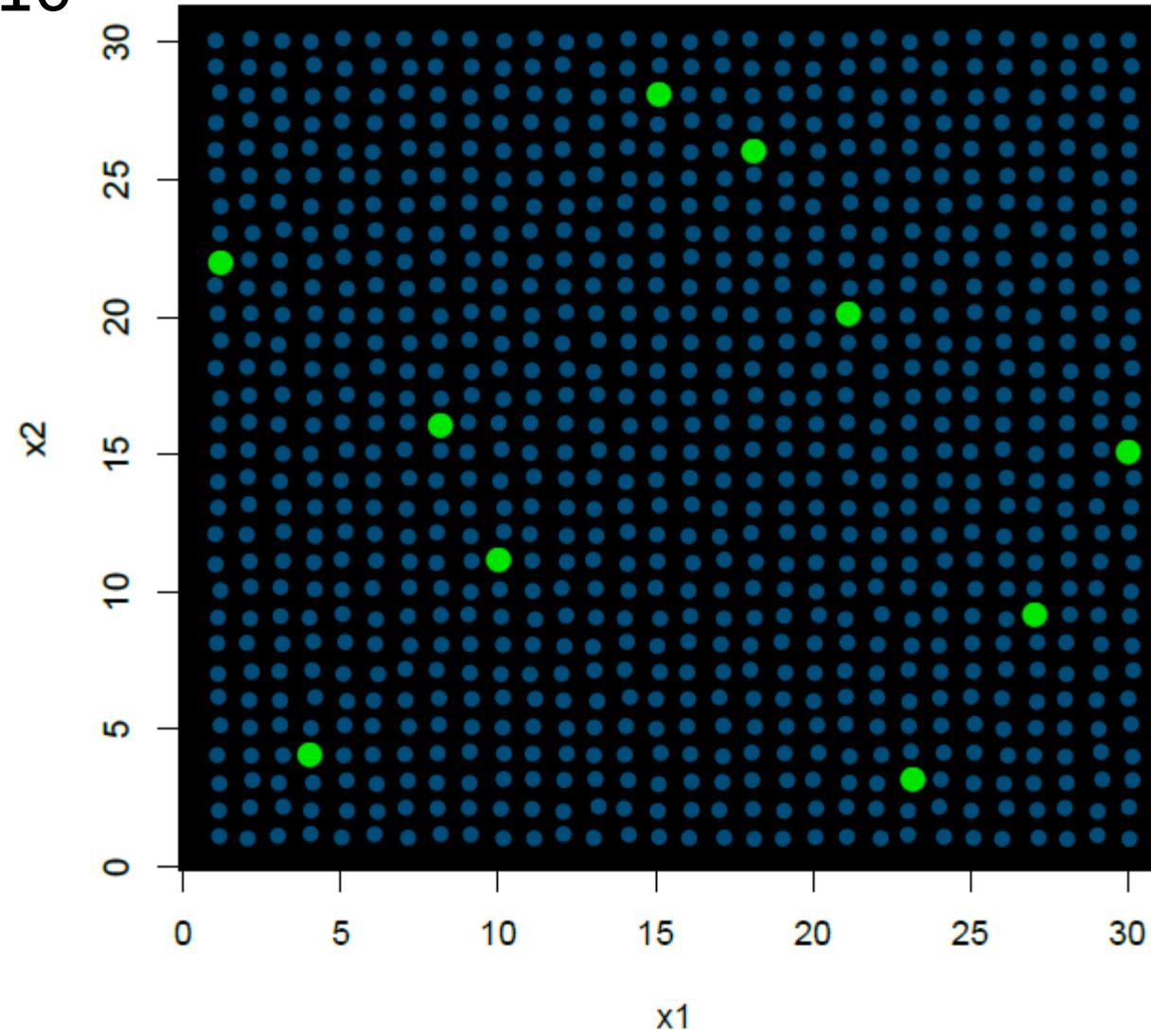


Calibration sampling

Sampling design

n = 10

contioned Latin Hypercube sampling (cLHS)

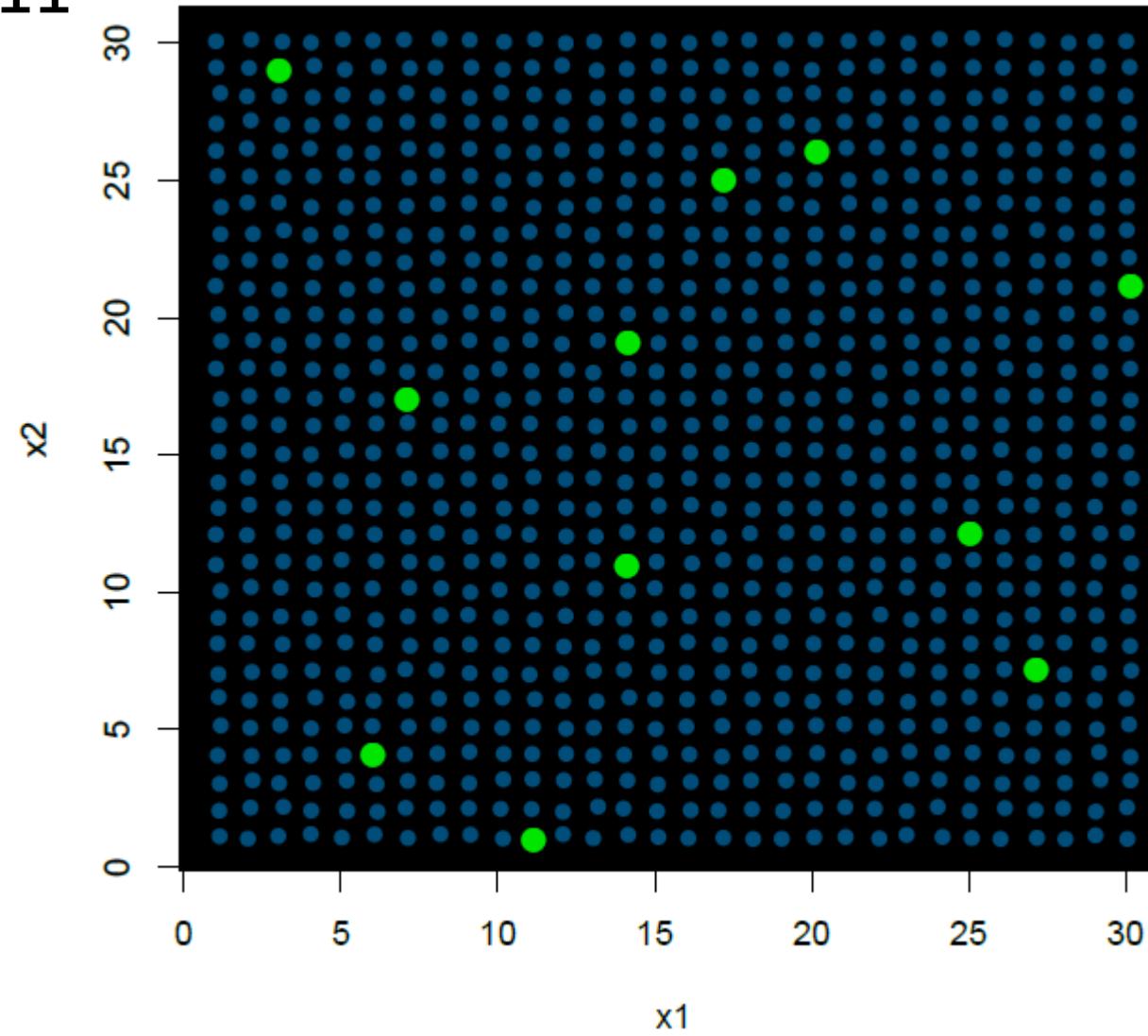


Calibration sampling

Sampling design

n = 11

contioned Latin Hypercube sampling (cLHS)

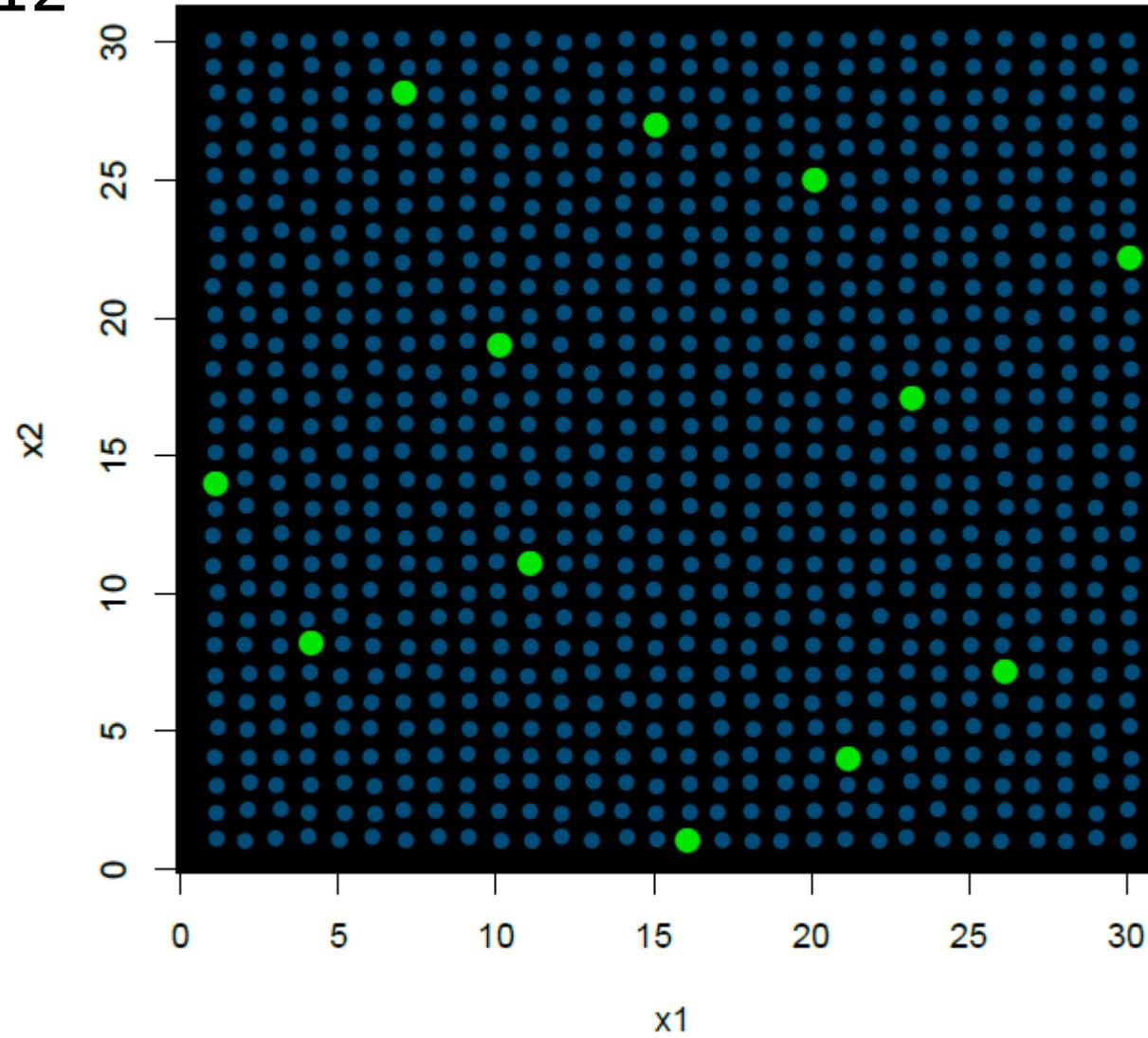


Calibration sampling

Sampling design

n = 12

contioned Latin Hypercube sampling (cLHS)

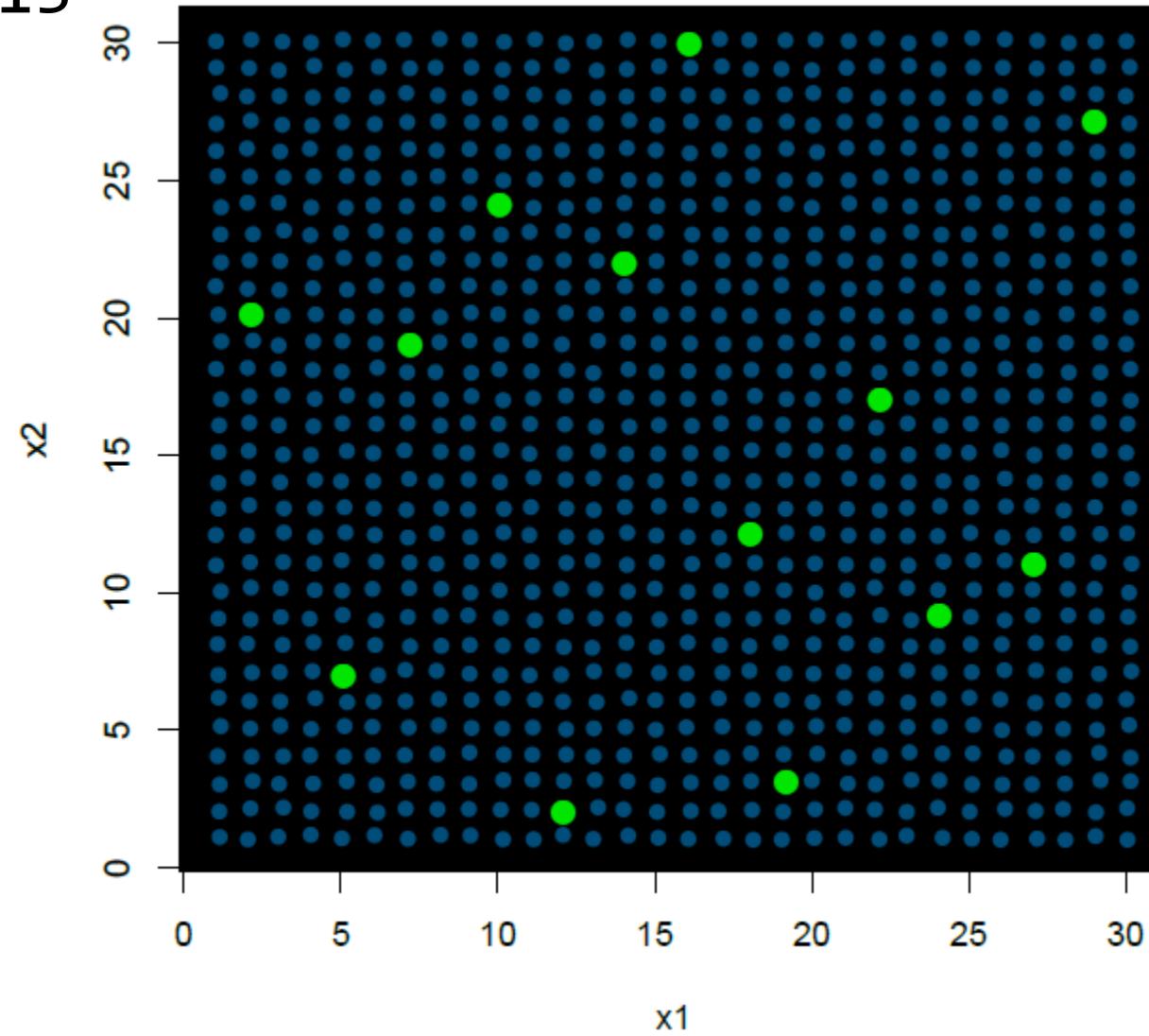


Calibration sampling

Sampling design

n = 13

contioned Latin Hypercube sampling (cLHS)

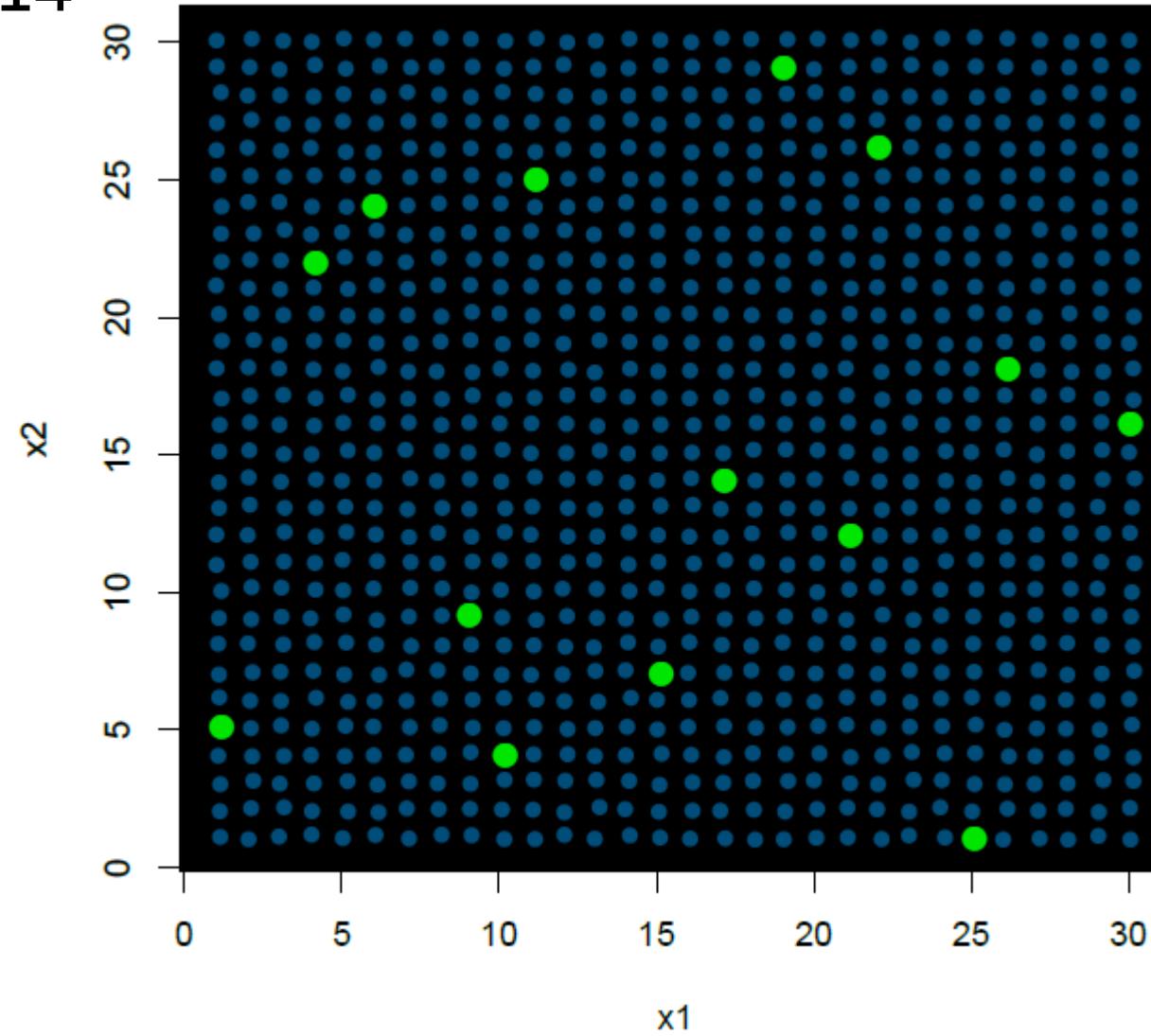


Calibration sampling

Sampling design

n = 14

contioned Latin Hypercube sampling (cLHS)

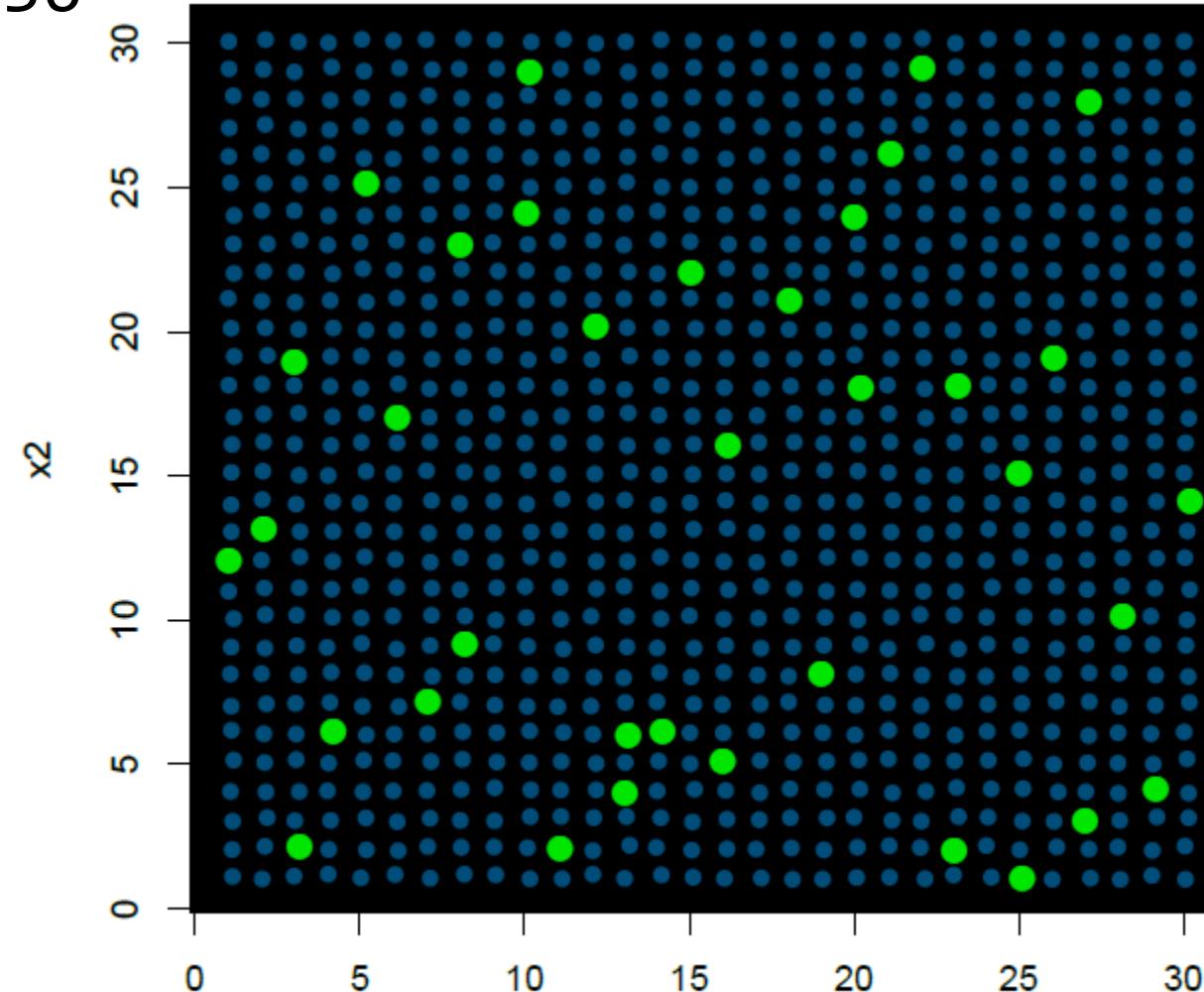


Calibration sampling

Sampling design

$n = 36$

contioned Latin Hypercube sampling (cLHS)

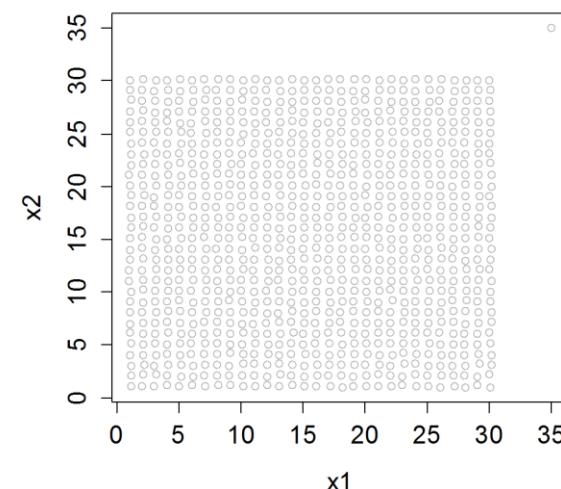
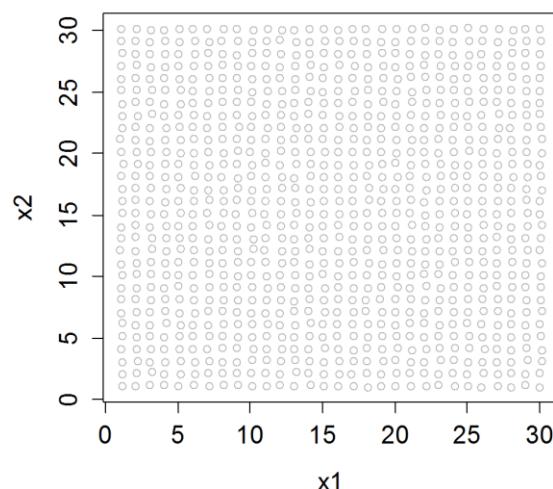


cLHS does not maximize the dissimilarity between samples

Calibration sampling

Sampling design

What if we have outliers in our data?

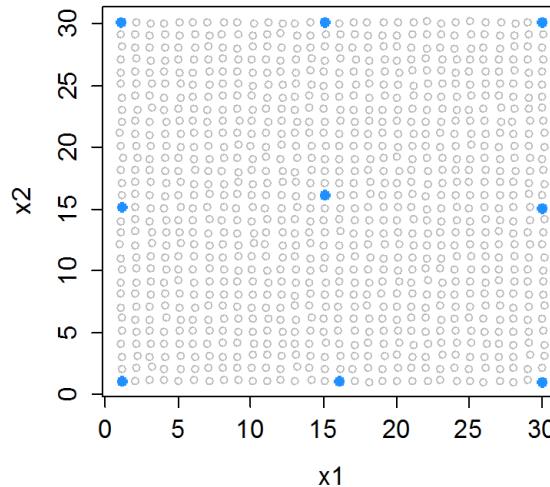


Calibration sampling

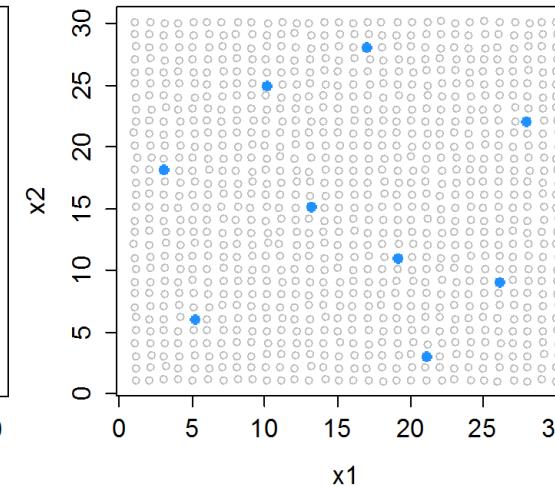
Sampling design

What if we have outliers in our data?

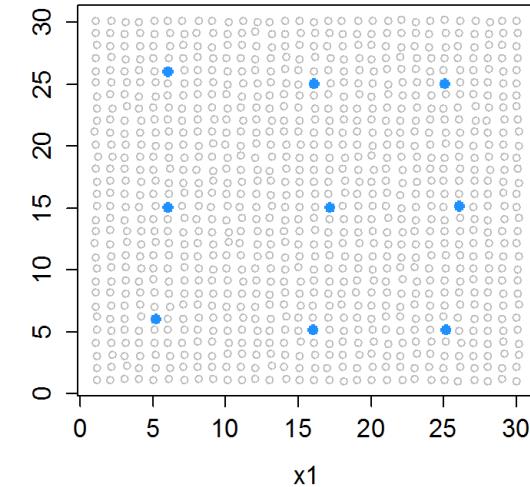
Kennard-Stone



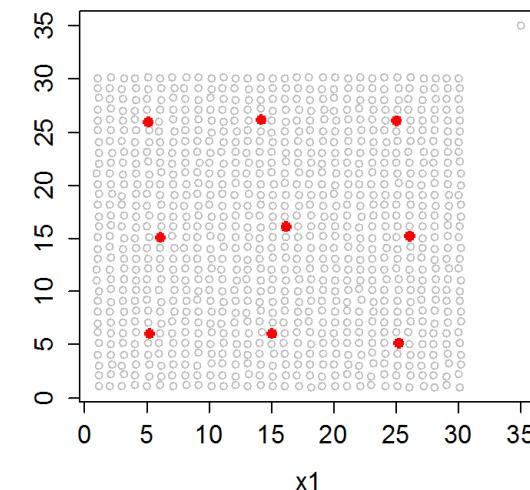
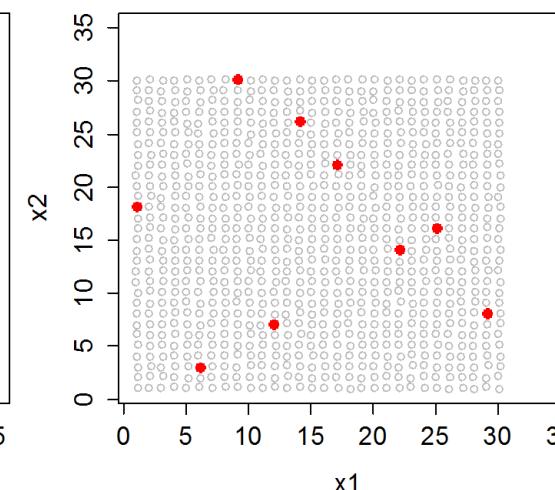
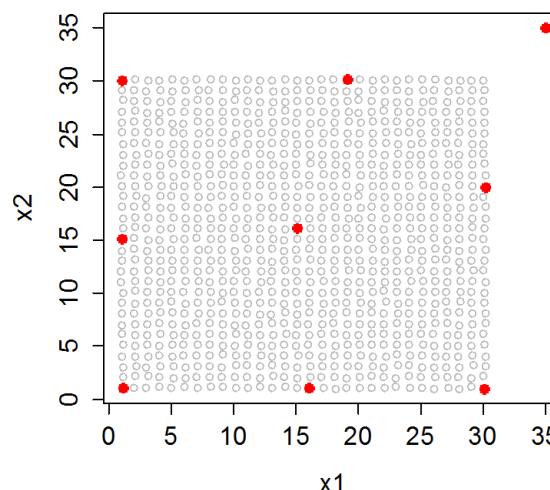
Latin hypercube



K-means



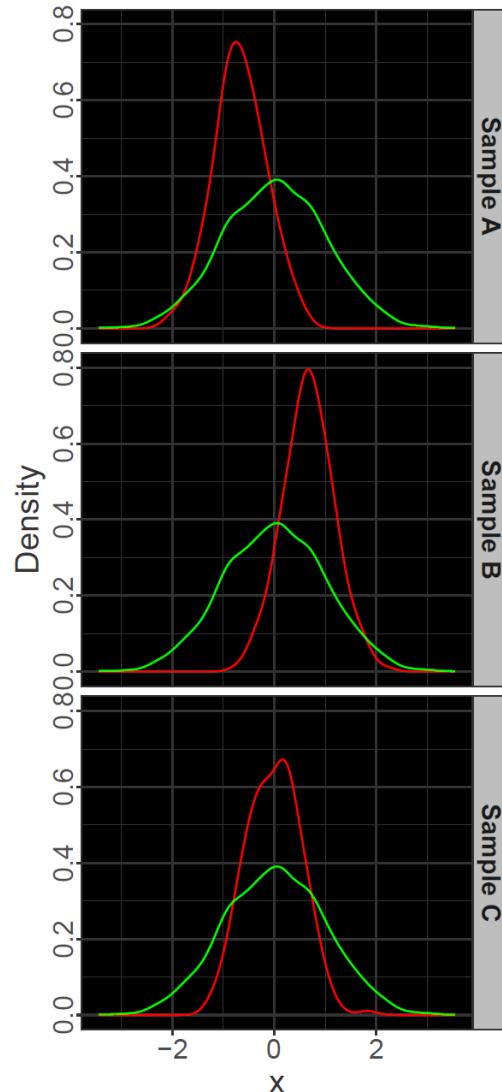
Without outliers



With one outlier

Calibration sampling

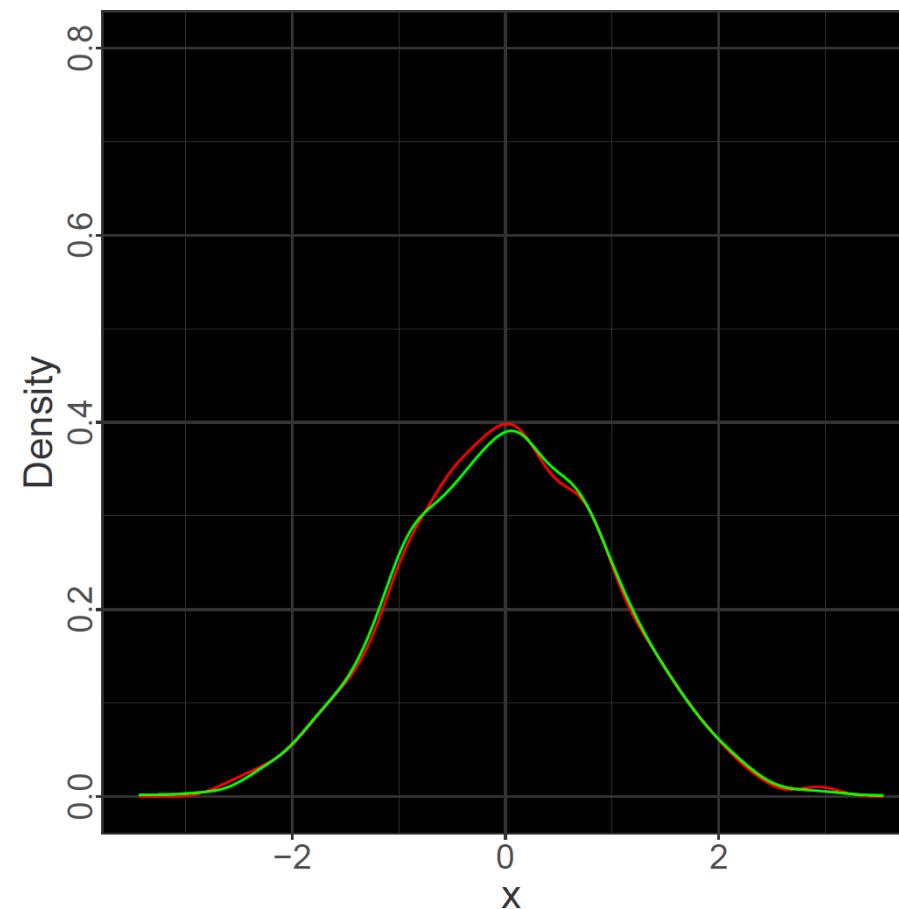
Low sampling efficiency



Size of the calibration set

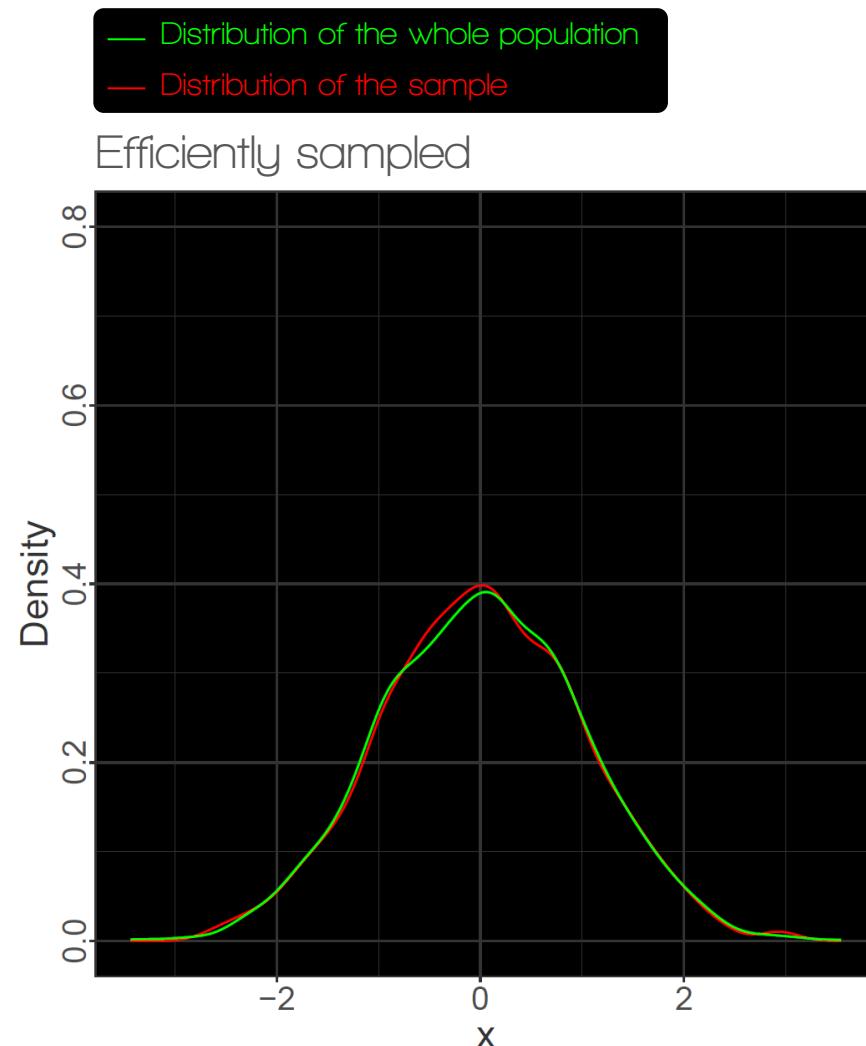
- Distribution of the whole population
- Distribution of the sample

Efficiently sampled



Calibration sampling

Size of the calibration set



Comparing the probability density function (pdf) of the population and the pdf of the sample set can be useful for assessing the representativeness of the sample set.

These pdfs can be computed for the PCs of the vis–NIR data.

Calibration sampling

The mean squared Euclidean distance (msd) can be used as a dissimilarity measure between the pdfs. It can be computed as follows:

$$msd = \frac{1}{k} \sum_{j=1}^k d^2 [P_s(x_j \in cs), P_p(x_j)]$$

where:

$$d^2 [P_s(x_j \in cs), P_p(x_j)] = \int_a^b [P_p(x_j) - P_s(x_j \in cs)]^2 dx_j$$

and where:

cs : a given subset of samples

$P_s(x_j \in cs)$: the estimated pdf of the j th PC of cs ,

$P_p(x_j)$: is the pdf of the j th PC for the whole population,

a and b : The range of the j th PC,

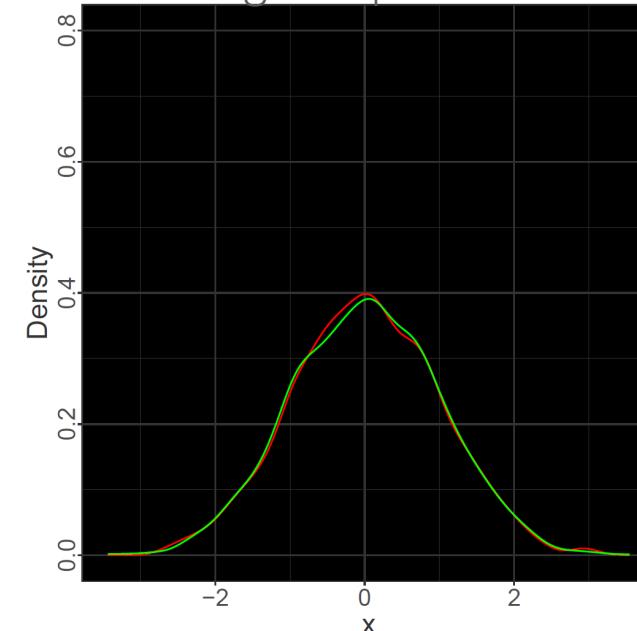
d^2 : represents the squared Euclidean distance between the two single distributions being compared,

k : the total number of PCs retained,

Size of the calibration set

— Distribution of the whole population
— Distribution of the sample

Efficiently sampled



Example...

Soil vis-NIR library of Europe ($n \sim 19.000$)

Calibration sampling algorithms

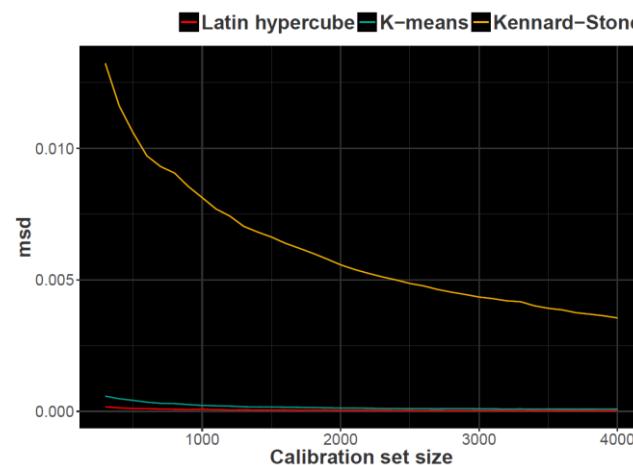
- Kennard-Stone sampling (KSS)
- Conditioned Latin hypercube sampling (cLHS)
- K-means sampling (KMS)

Calibration set size (CSS)

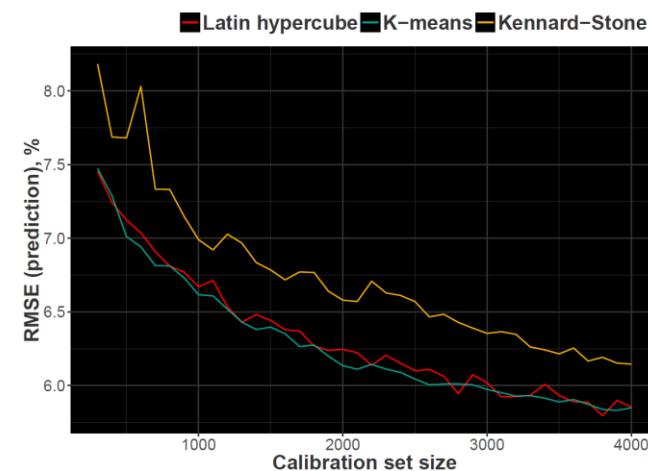
CSS = (300, 200, 400, ..., 4000)

The analysis is based on the comparison between the probability density function (pdf) of the population and the pdf of the sample set

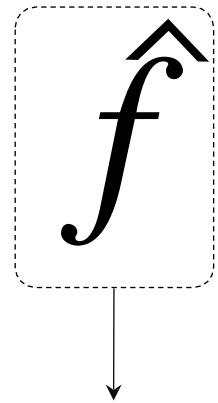
Comparisons of the pdfs



Prediction error



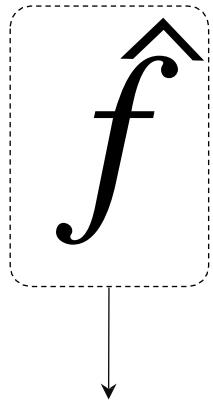
NIR Applications



Function that can be
approximated by using
chemometrics/machine learning

These functions are device-, product- and property-specific.

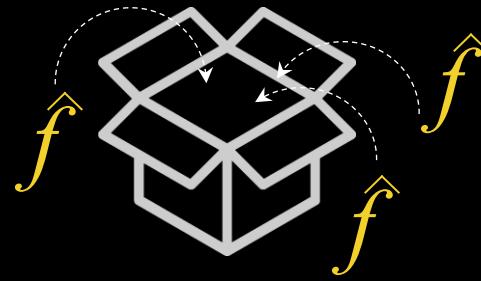
NIR Applications



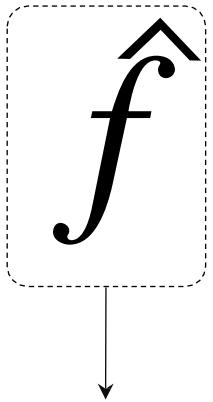
Function that can be approximated by using chemometrics/machine learning

These functions are device-, product- and property-specific.

For each product we build and group a set of functions that quantify relevant product properties (e.g. moisture, protein, fat). This group of functions is called application.



NIR Applications



Function that can be approximated by using chemometrics/machine learning

We currently use our own BUCHI software for that!

These functions are device-, product- and property-specific.

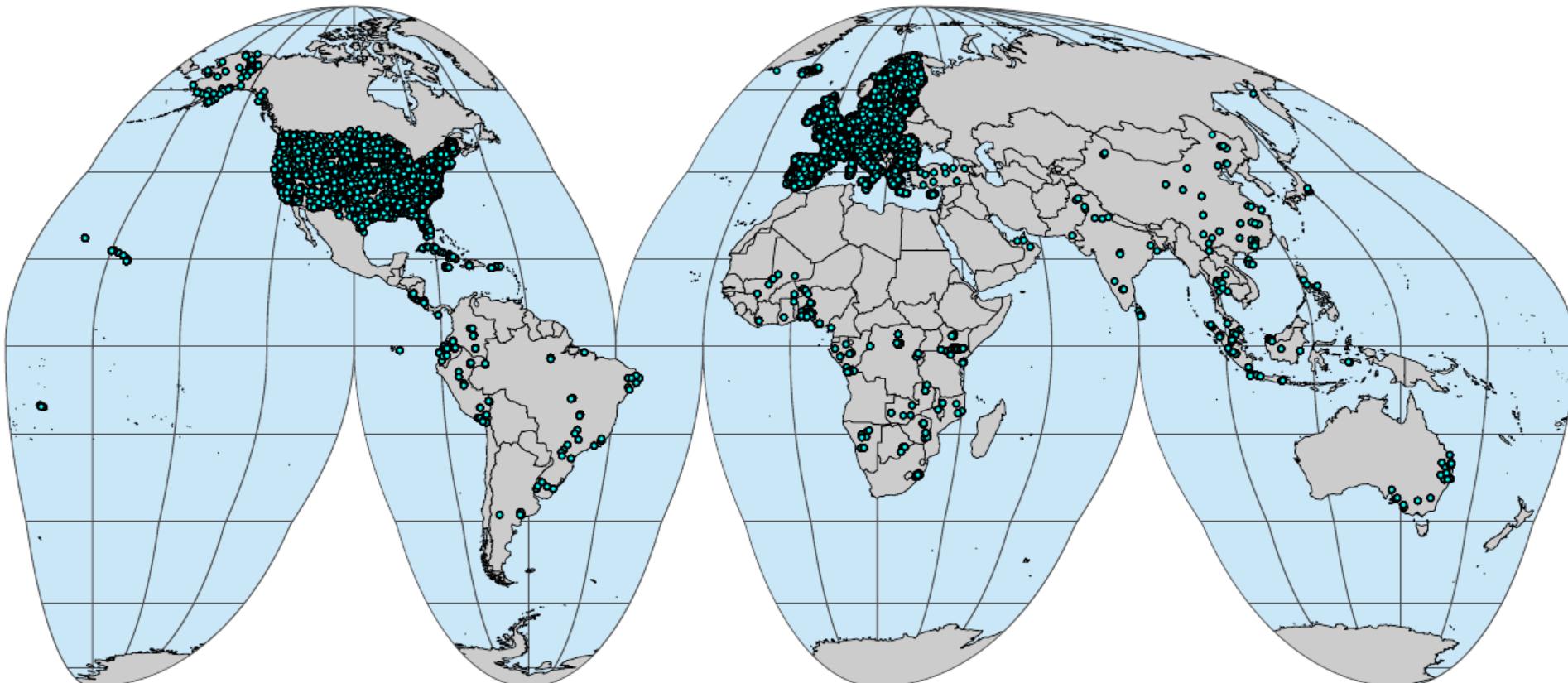
For each product we build and group a set of functions that quantify relevant product properties (e.g. moisture, protein, fat). This group of functions is called application.



From a user's perspective, a NIR device without an application can be perceived as useless.

How to benefit from spectral libraries to
make soil spectroscopy operational?

Dataset	Coverage	VNIR sample size
ICRAF_ISRIC	World	4438
KSSL	USA	19807
LUCAS	Europe	40818
Total		65063



Source: <https://soilspectroscopy.org/oss1-updates/>

We should aim at reducing entry barrier/implementation costs!!!

Implementation

It can be expensive

- Hardware and software costs
- Cost of reference analyses
- Elevated time investments
- People

Routine analysis

It can make boost the efficiency of your lab:

- Reduces need of conventional methods
- A sample is measured in matter of seconds
- Requires periodic validation and model/method updates

We should aim at reducing entry barrier/implementation costs!!!

Time to routine!

Implementation

It can be expensive

- Hardware and software costs
- Cost of reference analyses
- Elevated time investments
- People

Routine analysis

It can make boost the efficiency of your lab:

- Reduces need of conventional methods
- A sample is measured in matter of seconds
- Requires periodic validation and model/method updates

NIR/IR Spectroscopy – An equation for success

$$S_{\text{uccess}} = f(\text{hardware, software, data, models, people, automation, ...})$$

Contributions

- Focus on usability of spectral libraries rather than modeling
- Scale up soil spectroscopy
- Make spectroscopy simple and operational for "simple" domains (i.e. fast to implement and service).
- Think globally fit locally*

*Saul, L. K., & Roweis, S. T. (2003). Think globally, fit locally: unsupervised learning of low dimensional manifolds. *Journal of machine learning research*, 4(Jun), 119-155.

Success = *f* (hardware, software, data, chemometrics, automation, people, ...)

$S_{\text{uccess}} = f(\text{hardware, software, data, chemometrics, automation, people, ...})$

↓
Library search methods

Main users

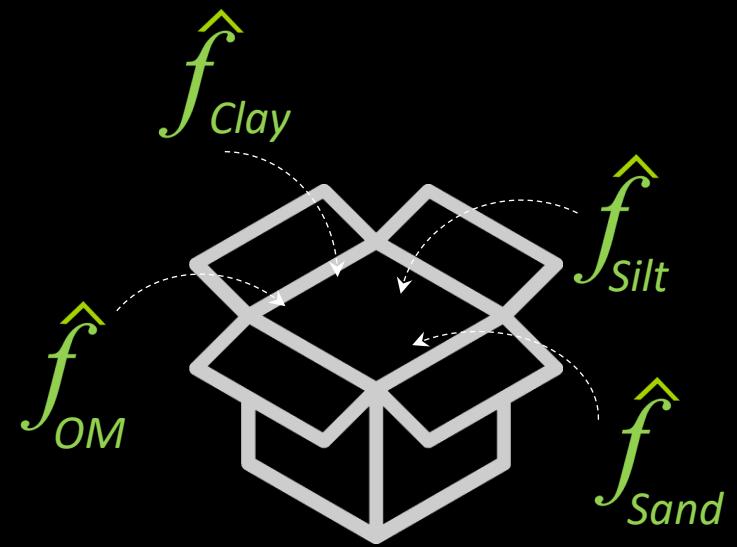
Big farms (soil, plant tissues, fertilizers, final product, etc)

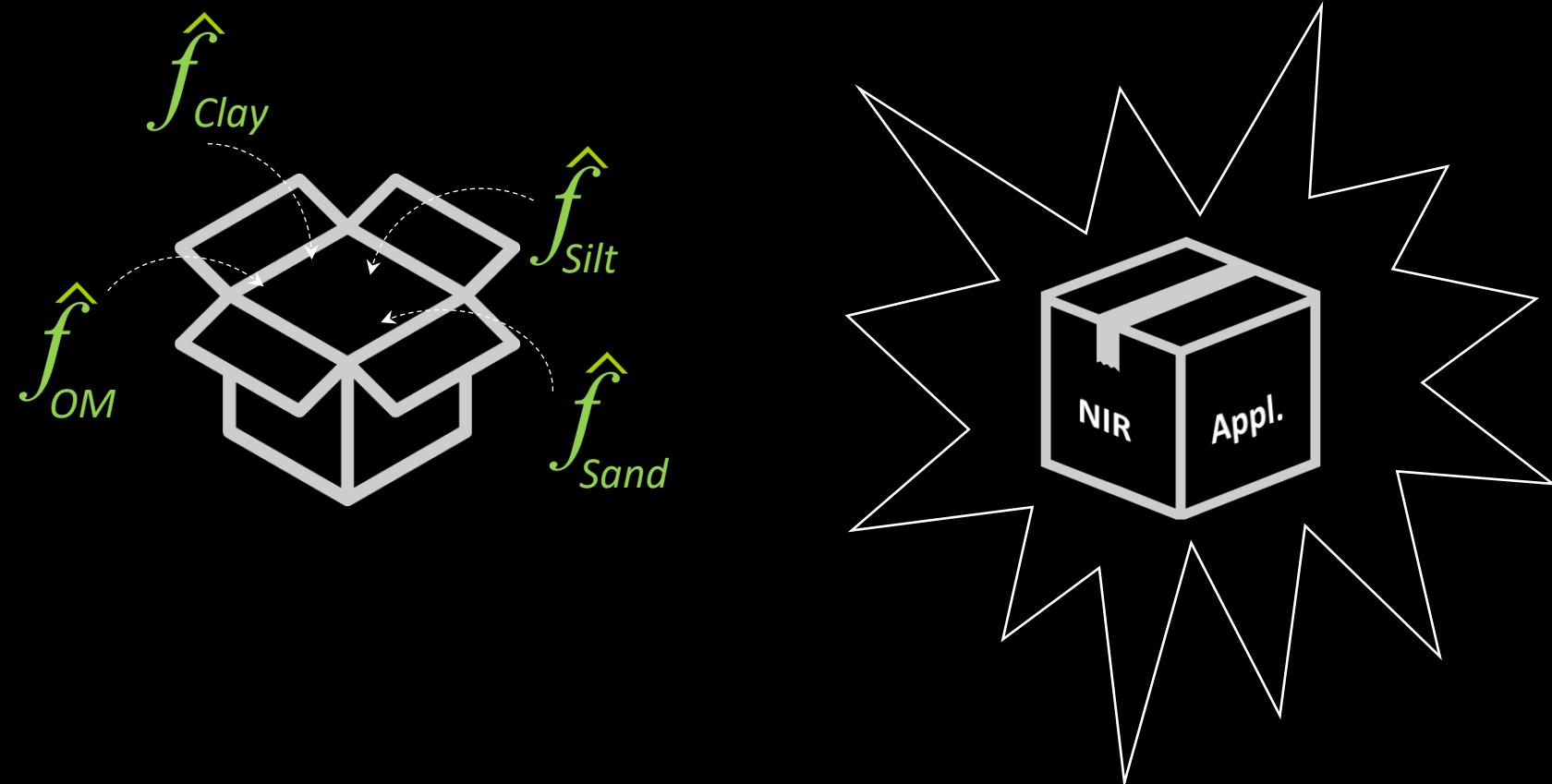
Service labs

Research institutes/Universities

NIR/IR Spectroscopy – An equation for success

$$S_{\text{uccess}} = f(\text{hardware, software, data, models, people, automation, ...})$$





debunking myths

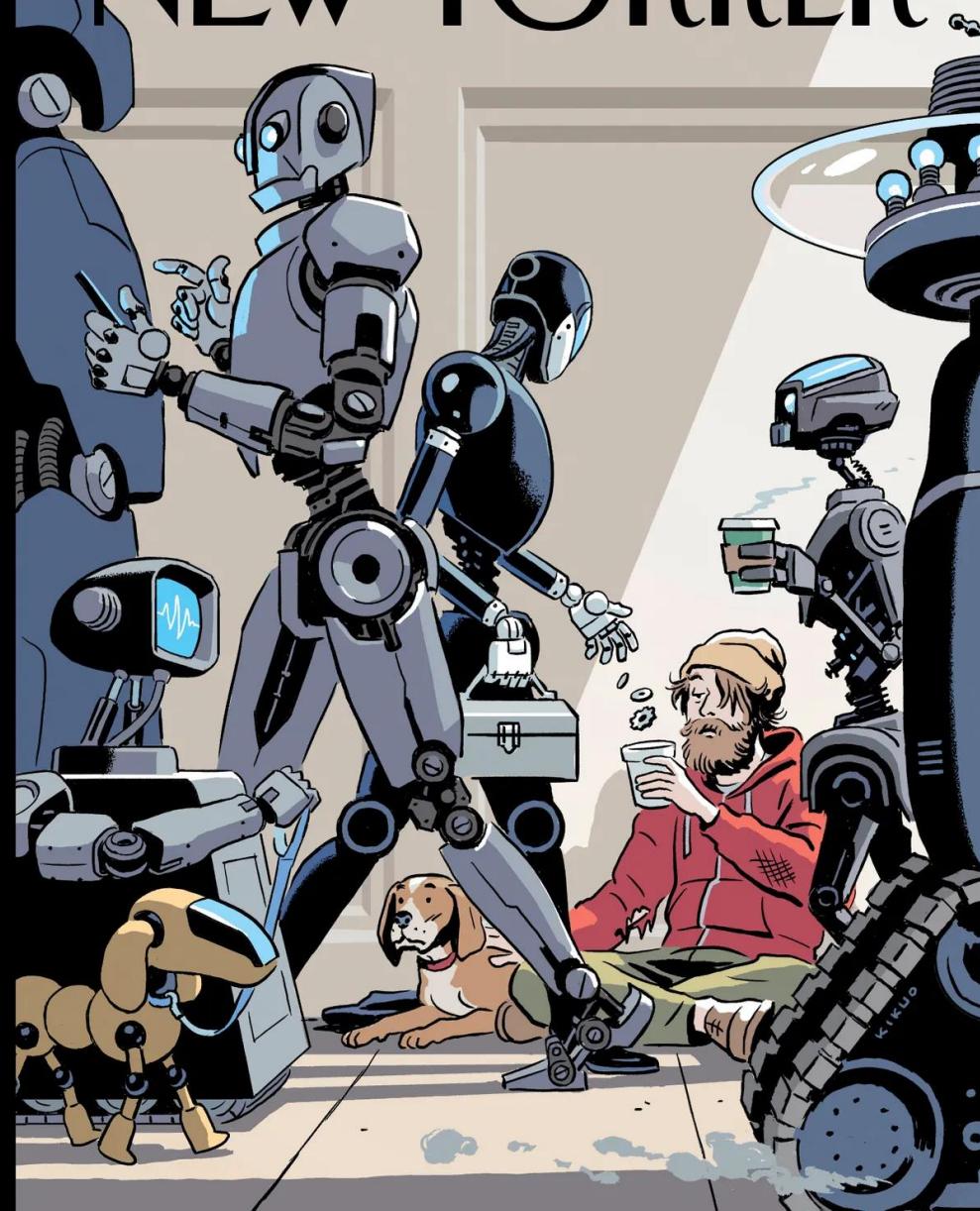
1. Could soil spectroscopy eventually **replace** conventional methods?

PRICE \$8.99

THE

OCT. 23, 2017

THE NEW YORKER





Purchase PDF

Find at Imperial

Access through another institution

Imperial College London does not subscribe to this content on ScienceDirect. X

Article preview

Abstract

Introduction

Section snippets

References (77)

Cited by (322)



ELSEVIER

Advances in Agronomy

Volume 132, 2015, Pages 139-159



Chapter Four - Soil Spectroscopy: An Alternative to Wet Chemistry for Soil Monitoring

Marco Nocita * , Antoine Stevens , Bas van Wesemael , Matt Aitkenhead , Martin Bachmann , Bernard Barthès , Eyal Ben Dor **, David J. Brown , Michael Clairotte , Adam Csorba , Pierre Dardenne , Jose A.M. Demattè ##, Valerie Genot , Cesar Guerrero ***, Maria Knadel , Luca Montanarella *, Carole Noon , Leonardo Ramirez-Lopez , Jean Robertson , Hiro Sakai , Johanna Wetterling ...

Show more ▾

Recommended articles ^

[Occurrence, Detection, and Molecular and Metabolic Characterization of...](#)

Advances in Agronomy, Volume 132, 2015, pp. 1...
Magdalena Frqc, ..., Takashi Yaguchi

[Hydrological Aspects of Arsenic Contamination of Groundwater in...](#)

Advances in Agronomy, Volume 132, 2015, pp. 7...
Saugata Datta

[Do we really need large spectral libraries for local scale SOC...](#)

Soil and Tillage Research, Volume 155, 2016, p...
César Guerrero, ..., Raphael A. Viscarra Rossel

Show 3 more articles ▾



Purchase PDF

Find at Imperial

Access through another institution

ⓘ Imperial College London does not subscribe to this content on ScienceDirect.

X

Article preview

Abstract

Introduction

Section snippets

References (77)

Cited by (322)



Chapter Four Spectroscopy: An Alternative to Electrochemistry for Soil Monitoring

Marco Nocita * §, Antoine Stevens §, Bas van Wesemael §, Matt Aitkenhead ¶, Martin Bachmann ||, Bertrand Barthès #, Eyal Ben Dor **, David J. Brown §§, Michael Clairotte #, Adam Csorba ¶¶, Pierre Dardenne |||, Jose A.M. Demattè §§, Valerie Genot †, Cesar Guerrero ***, Maria Knadel §§§, Luca Montanarella *, Carole Noon §, Leonardo Ramirez-Lopez ¶¶¶, Jean Robertson ¶, Hiro Sakai ||||..., Johanna Wetterling §§§§

Show more ▾

Recommended articles

[Occurrence, Detection, and Molecular and Metabolic Characterization of...](#)

Advances in Agronomy, Volume 132, 2015, pp. 1...
Magdalena Frqc, ..., Takashi Yaguchi

[Hydrological Aspects of Arsenic Contamination of Groundwater in...](#)

Advances in Agronomy, Volume 132, 2015, pp. 7...
Saugata Datta

[Do we really need large spectral libraries for local scale SOC...](#)

Soil and Tillage Research, Volume 155, 2016, p...
César Guerrero, ..., Raphael A. Viscarra Rossel

Show 3 more articles ▾

Example: NIR spectroscopy for soil mapping and management

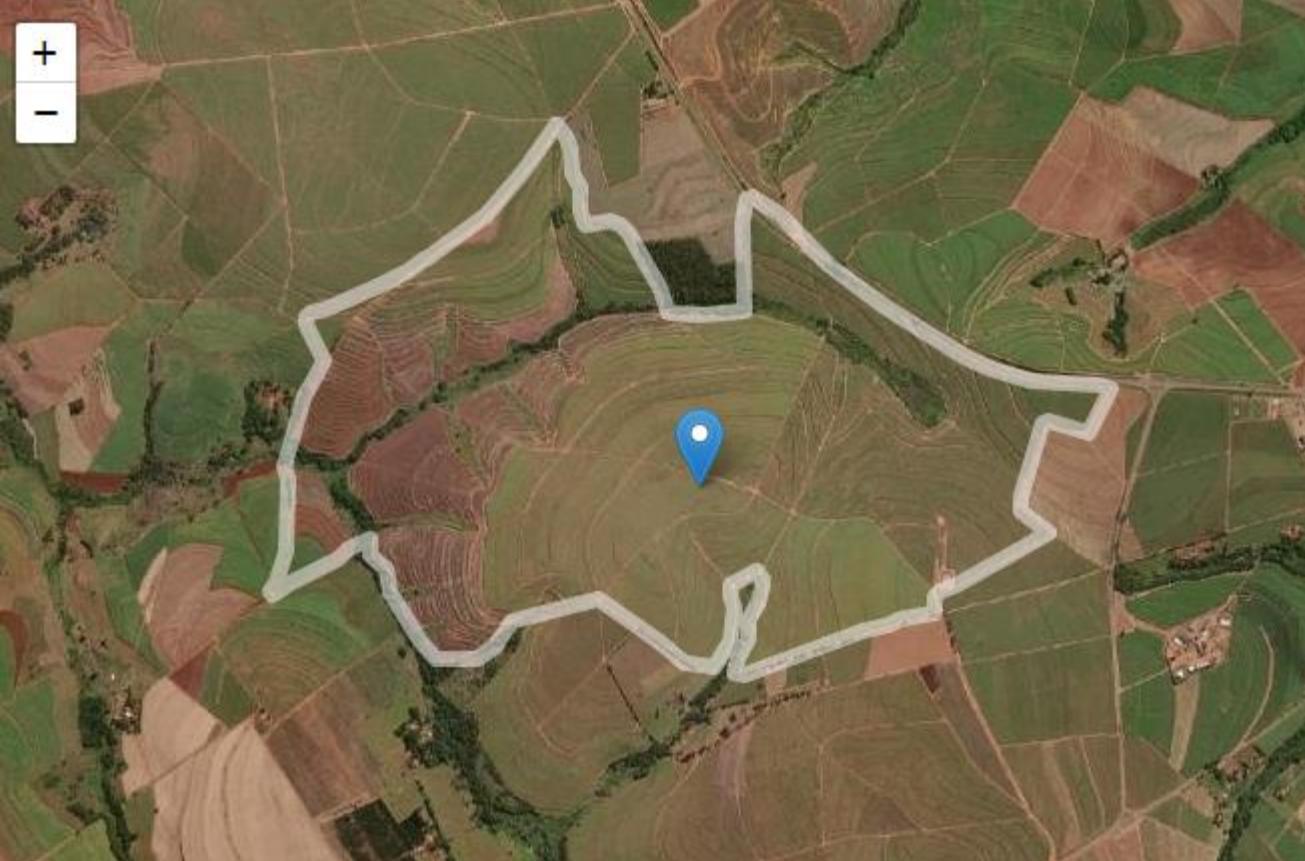
Ramirez-Lopez, Leonardo, et al. "Robust soil mapping at the farm scale with vis-NIR spectroscopy." *European Journal of Soil Science* 70.2 (2019): 378-393.

Example:

473 ha in the municipality of Barra Bonita (Brazil)

Sampling at two depths (A: 0 – 20 cm; B: 80 – 100 cm)

Total samples: 910



Example:

473 ha in the municipality of Barra Bonita (Brazil)

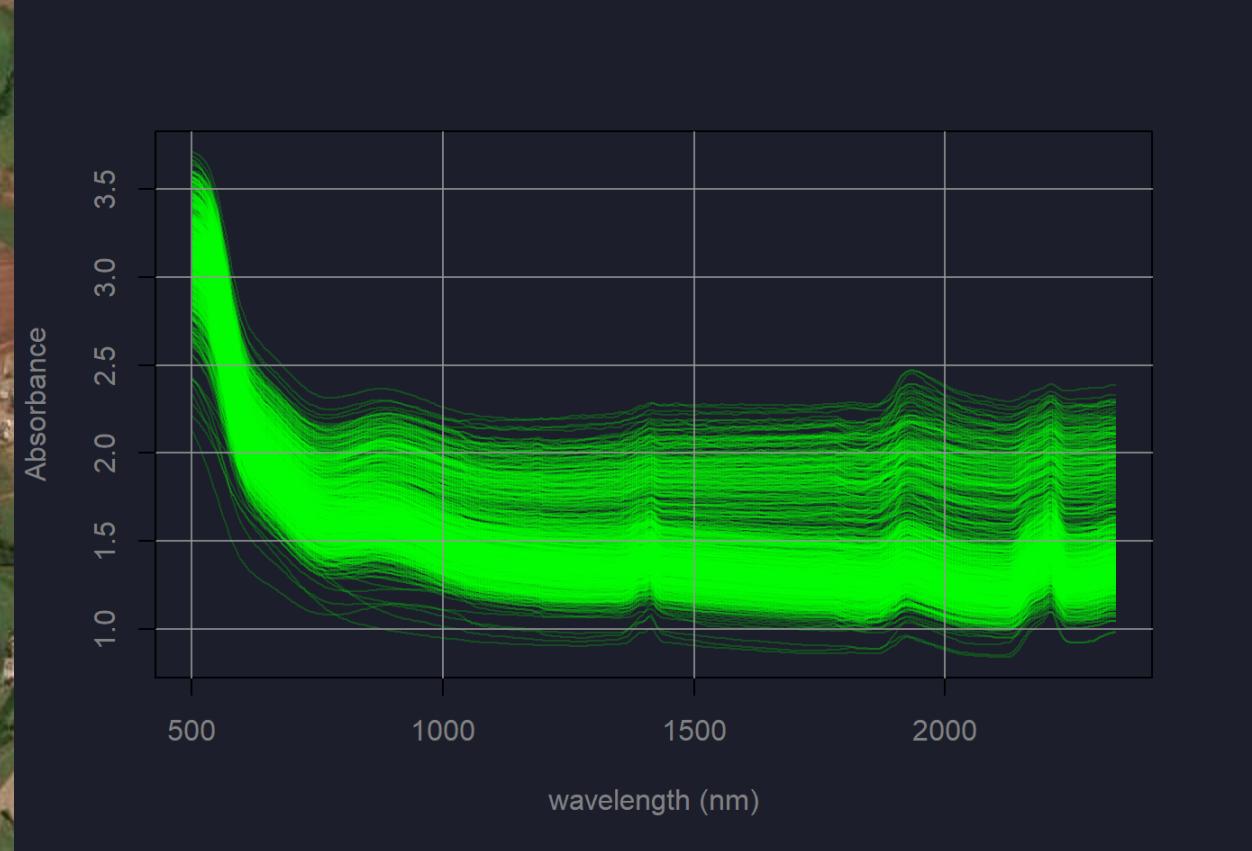
Sampling at two depths (A: 0 – 20 cm; B: 80 – 100 cm)

Total samples: 910



Sampling at two depths (A: 0 – 20 cm; B: 80 – 100 cm)

Total samples: 910



Example:

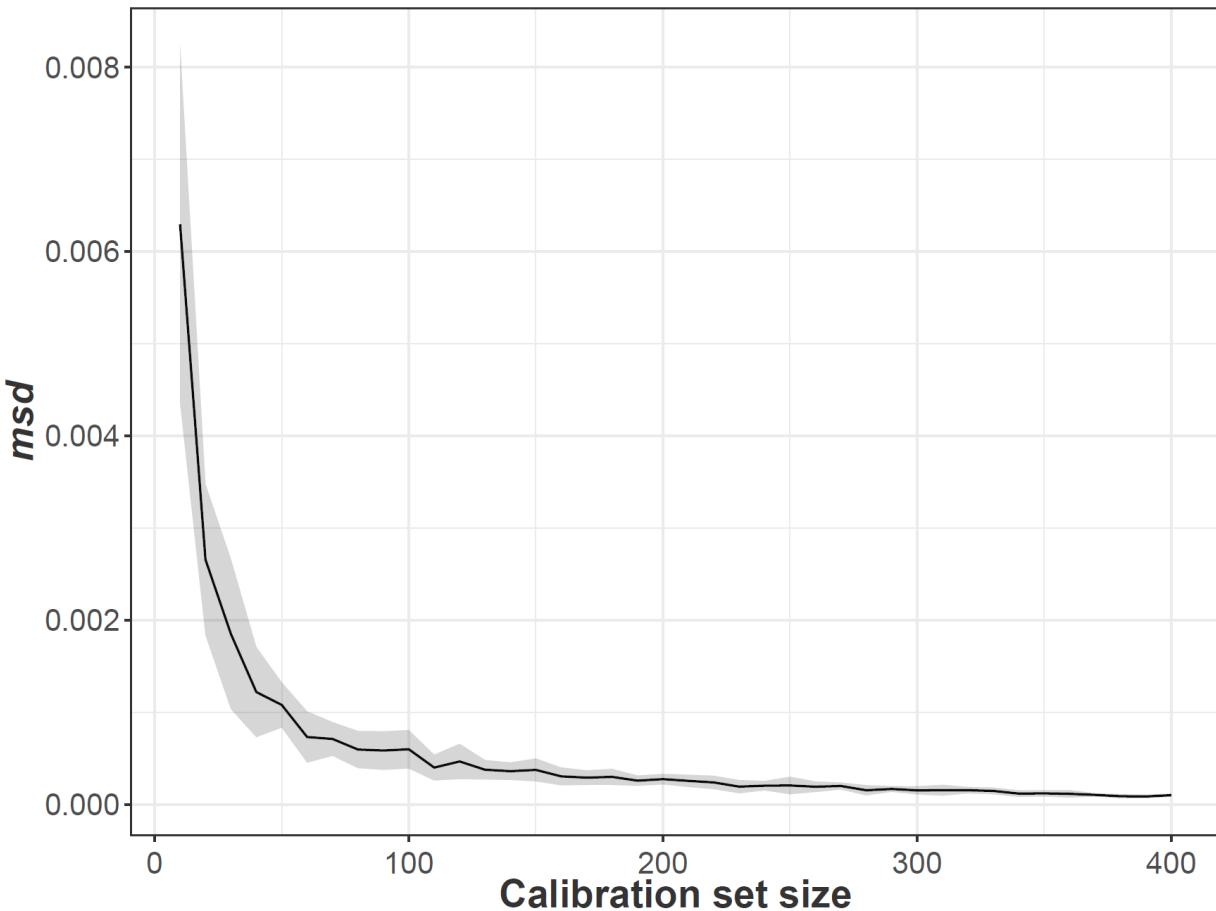
473 ha in the municipality of Barra Bonita (Brazil)

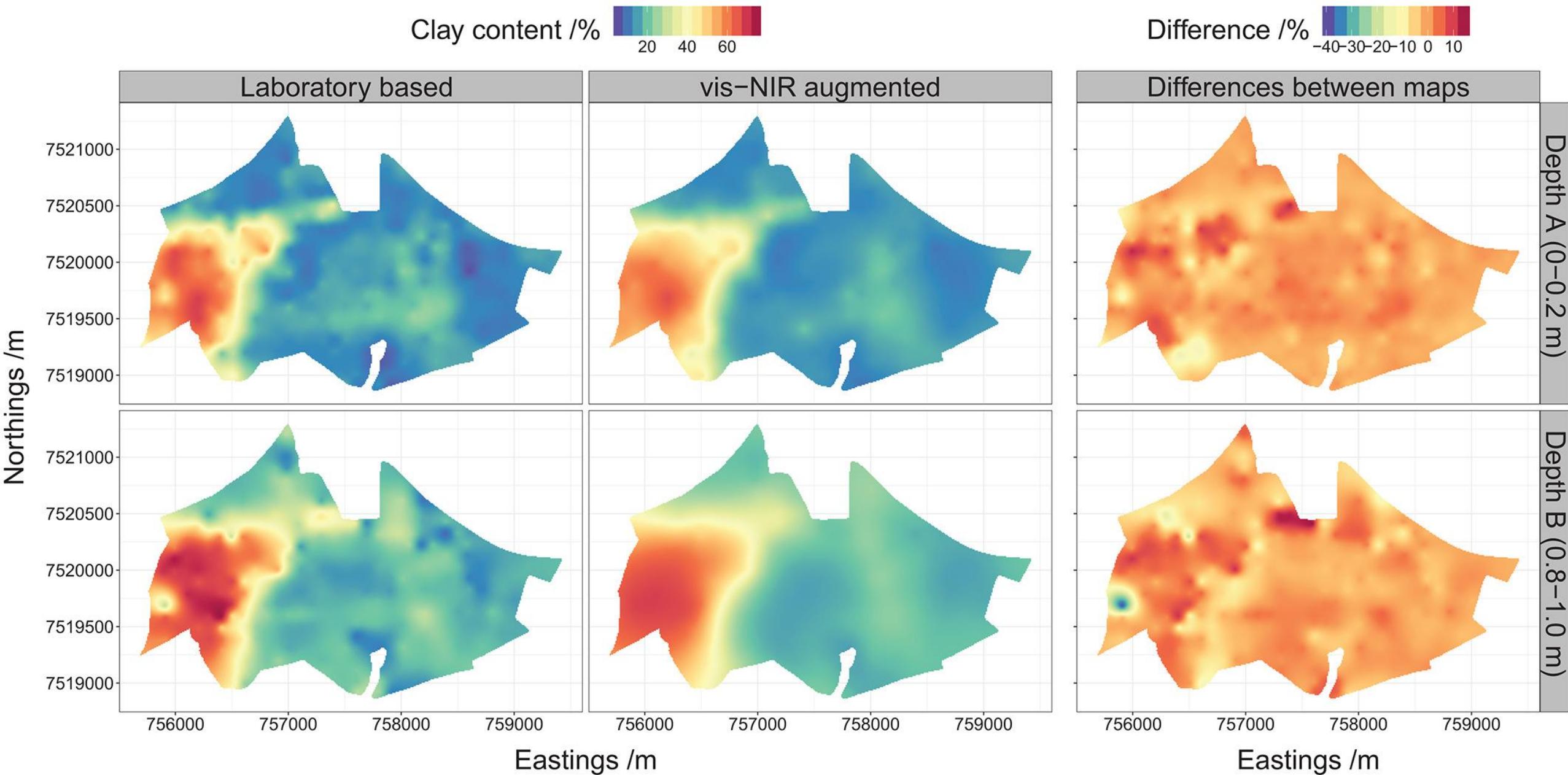
Sampling at two depths (A: 0 – 20 cm; B: 80 – 100 cm)

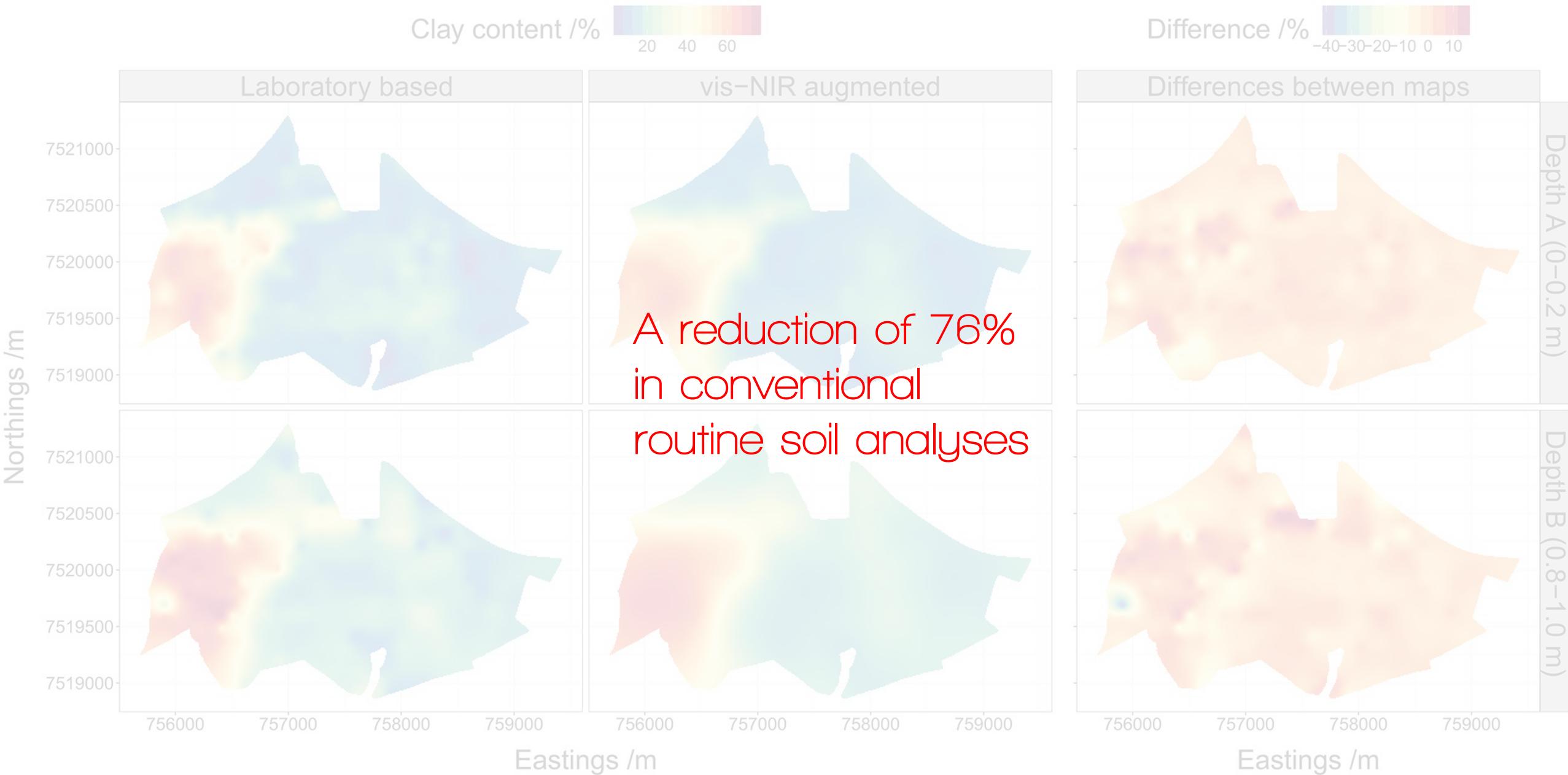
Total samples: 910



Optimal calibration size: 180 samples







Scaling soil spectroscopy

Short term
Learn and adjust

Scaling soil spectroscopy

Short term

Learn and adjust

Mid term

Move to modest amount of samples
processed with NIR/IR

Scaling soil spectroscopy

Short term

Learn and adjust

Mid term

Move to modest amount of samples
processed with NIR/IR

Long term

Make it fully operational for your
properties and move to a significant
amount of samples processed with
NIR/IR (get all your investments back)

Soil spectroscopy **complements** conventional methods (to different extents)

2. Is soil spectroscopy **cheap** and **fast**?

Mmm

Mmmmmmmmmmmmmmm

Mmmmmmmmmmmmmmm,
Yeah BUT

Mmmmmmmmmmmmmmm,
Yeah BUT not really

Mmmmmmmmmmmmmmm,
Yeah BUT not really

Implementation

Routine analysis

Mmmmmmmmmmmmmmm,
Yeah BUT not really

Implementation

It can be expensive

- Hardware and software costs
- Cost of reference analyses
- Elevated time investments
- People

Routine analysis

Mmmmmmmmmmmmmmm,
Yeah BUT not really

Implementation

It can be expensive

- Hardware and software costs
- Cost of reference analyses
- Elevated time investments
- People

Routine analysis

It can make boost the efficiency of your lab:

- Reduces need of conventional methods
- A sample is measured in matter of seconds
- Requires periodic validation and model/method updates

3. Can all important soil properties be quantified?

Let's ask a friend who has read TONS of scientific literature

Sound scientific evidence
[consistently showing good results]

Under research

Sound scientific evidence
[consistently showing good results]

Examples:

- Particle size
- Organic C
- N

Under research

Where do you want go?

Sound scientific evidence
[consistently showing good results]

Under research

4. Are global/universal soil spectroscopy models feasible?

Svalbard, Norway



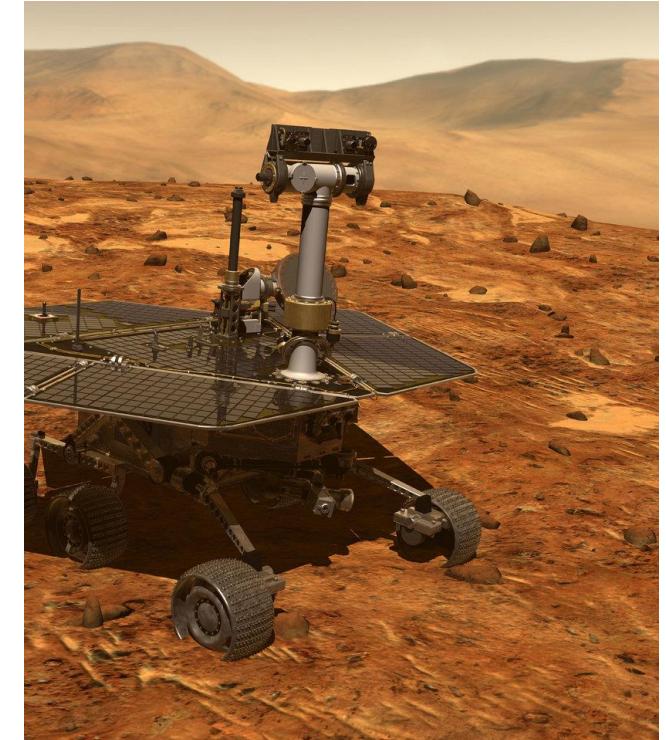
Svalbard, Norway



Uganda



Mars



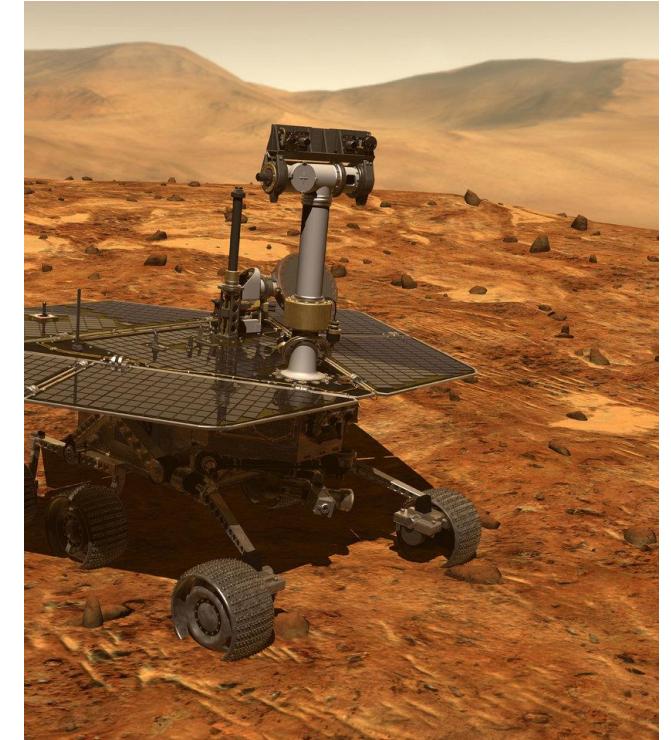
Svalbard, Norway



Uganda



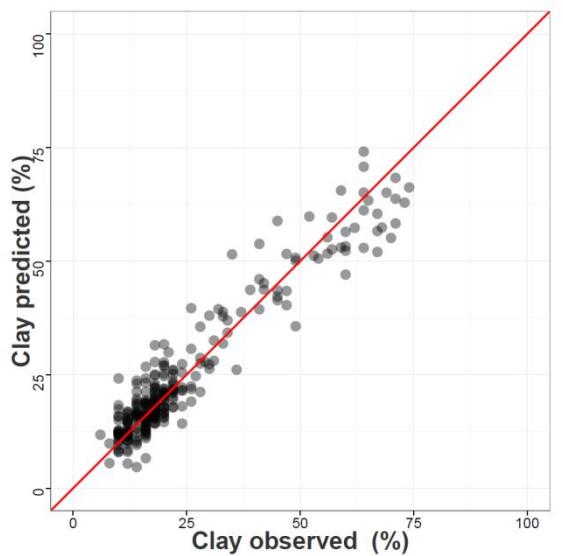
Mars



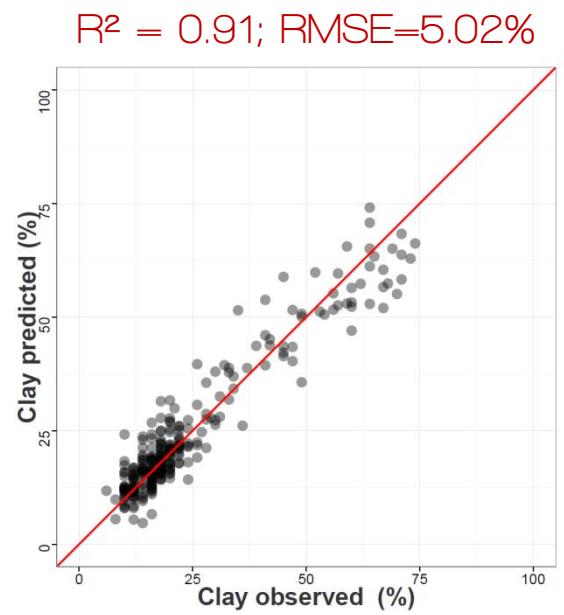
➔ Different soils might require different empirical models

Field scale

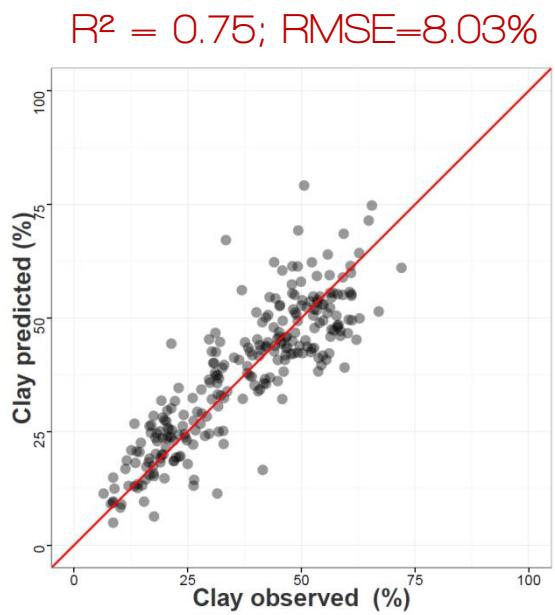
$R^2 = 0.91$; RMSE=5.02%



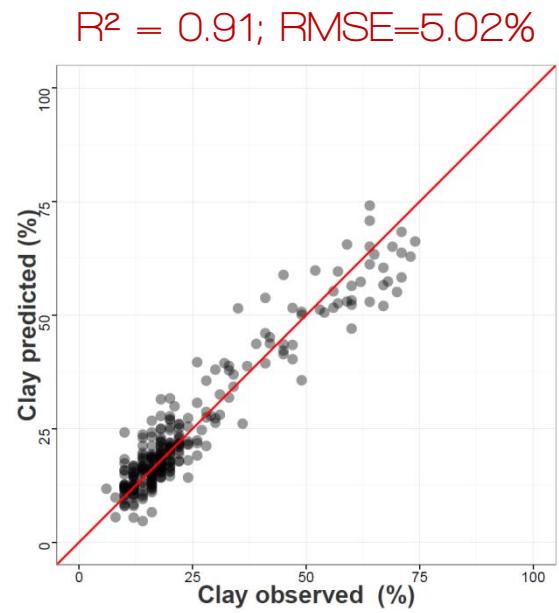
Field scale



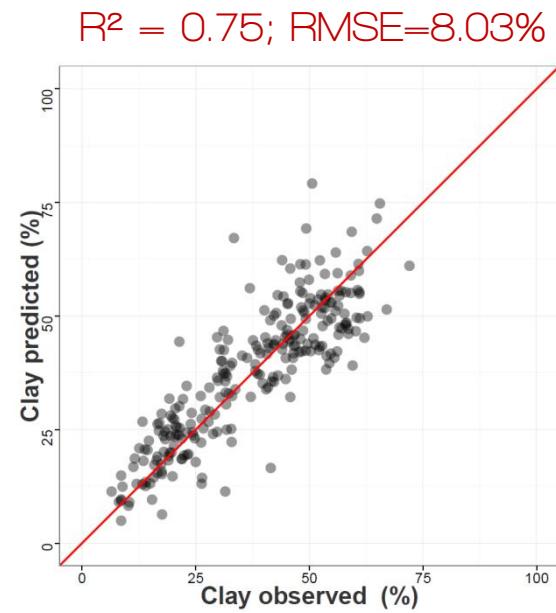
Regional scale



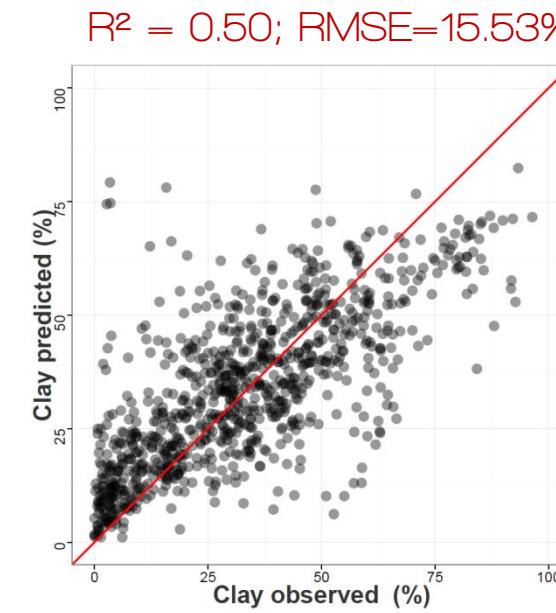
Field scale



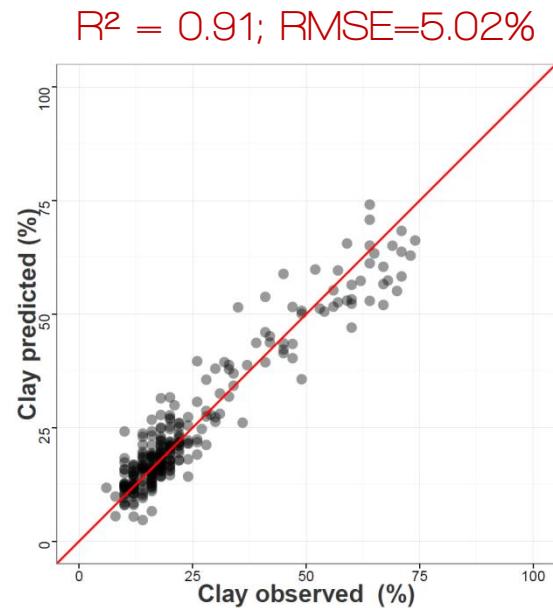
Regional scale



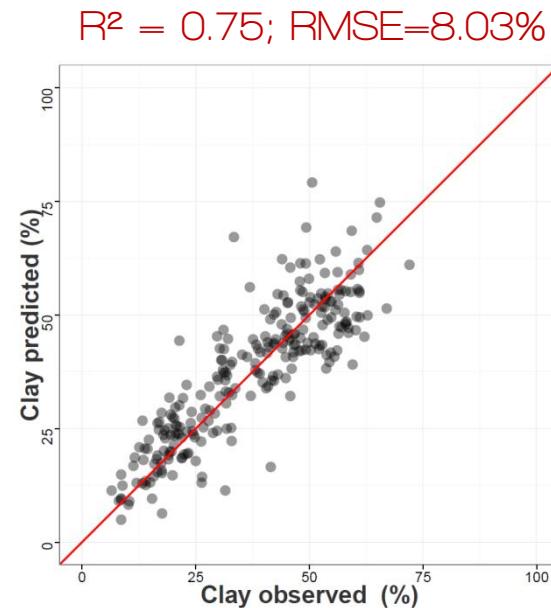
Global scale



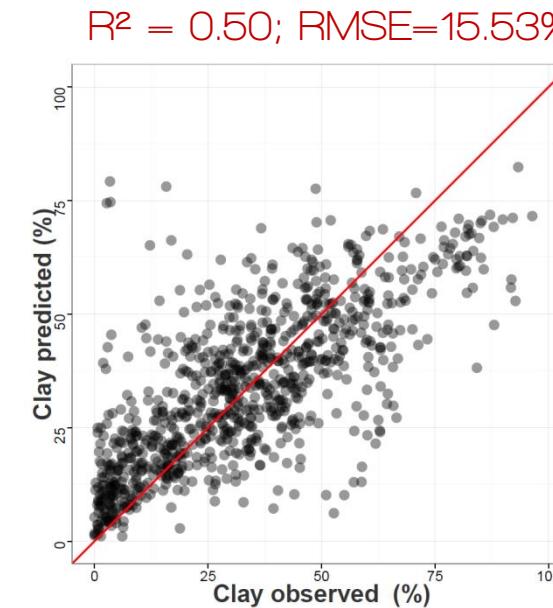
Field scale



Regional scale



Global scale

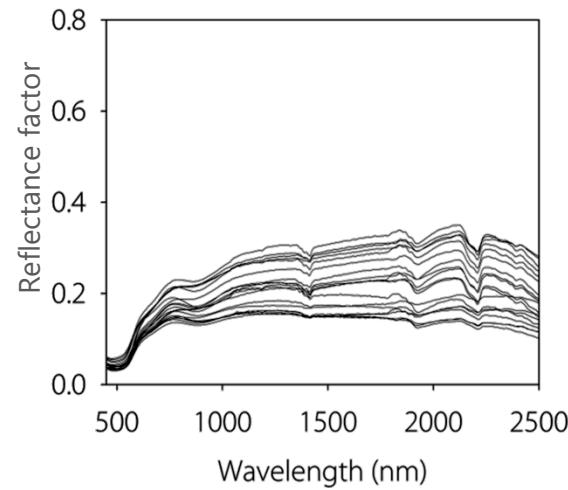


Data complexity

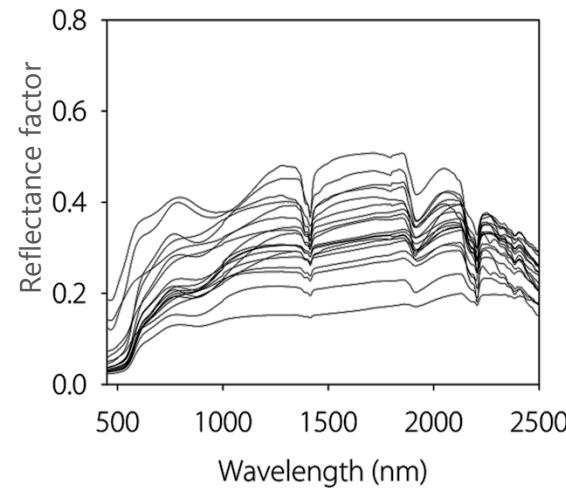


Degradation of the accuracy

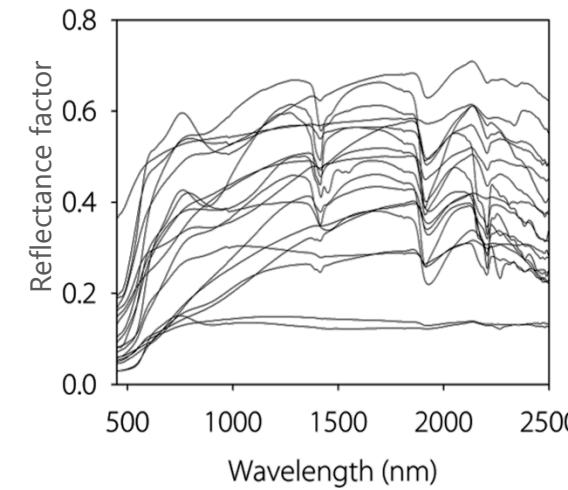
Field scale



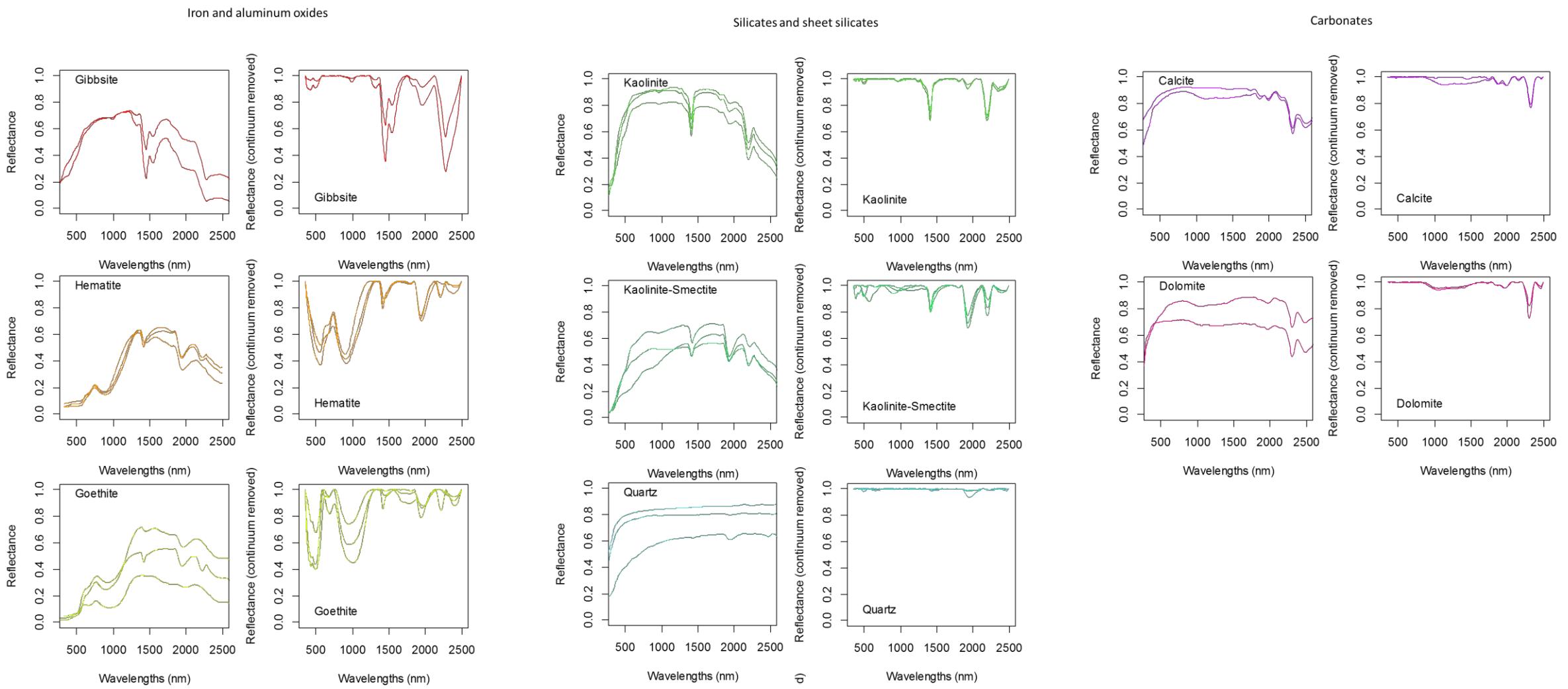
Regional scale

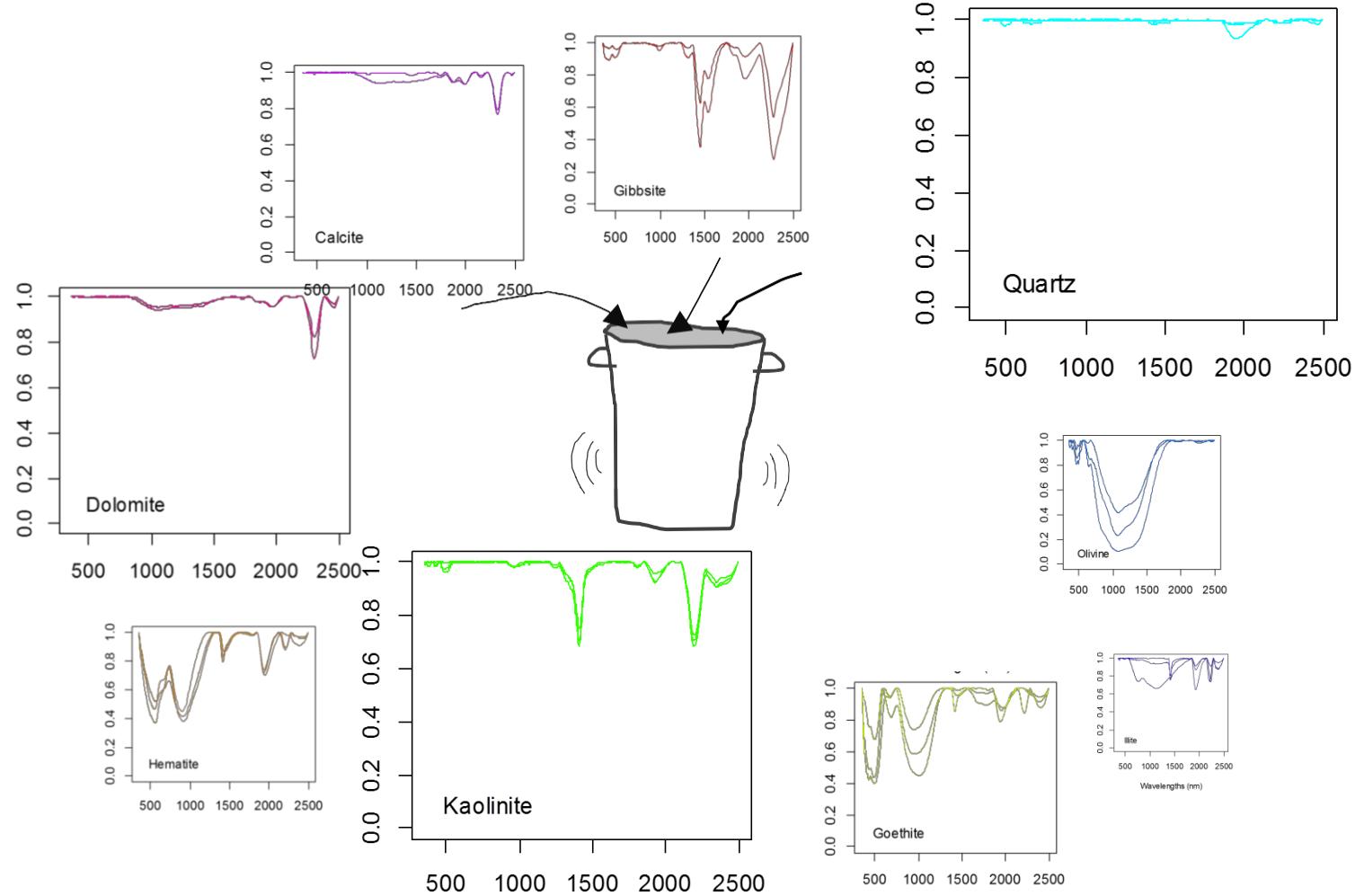


Global scale



20 spectra sampled at random



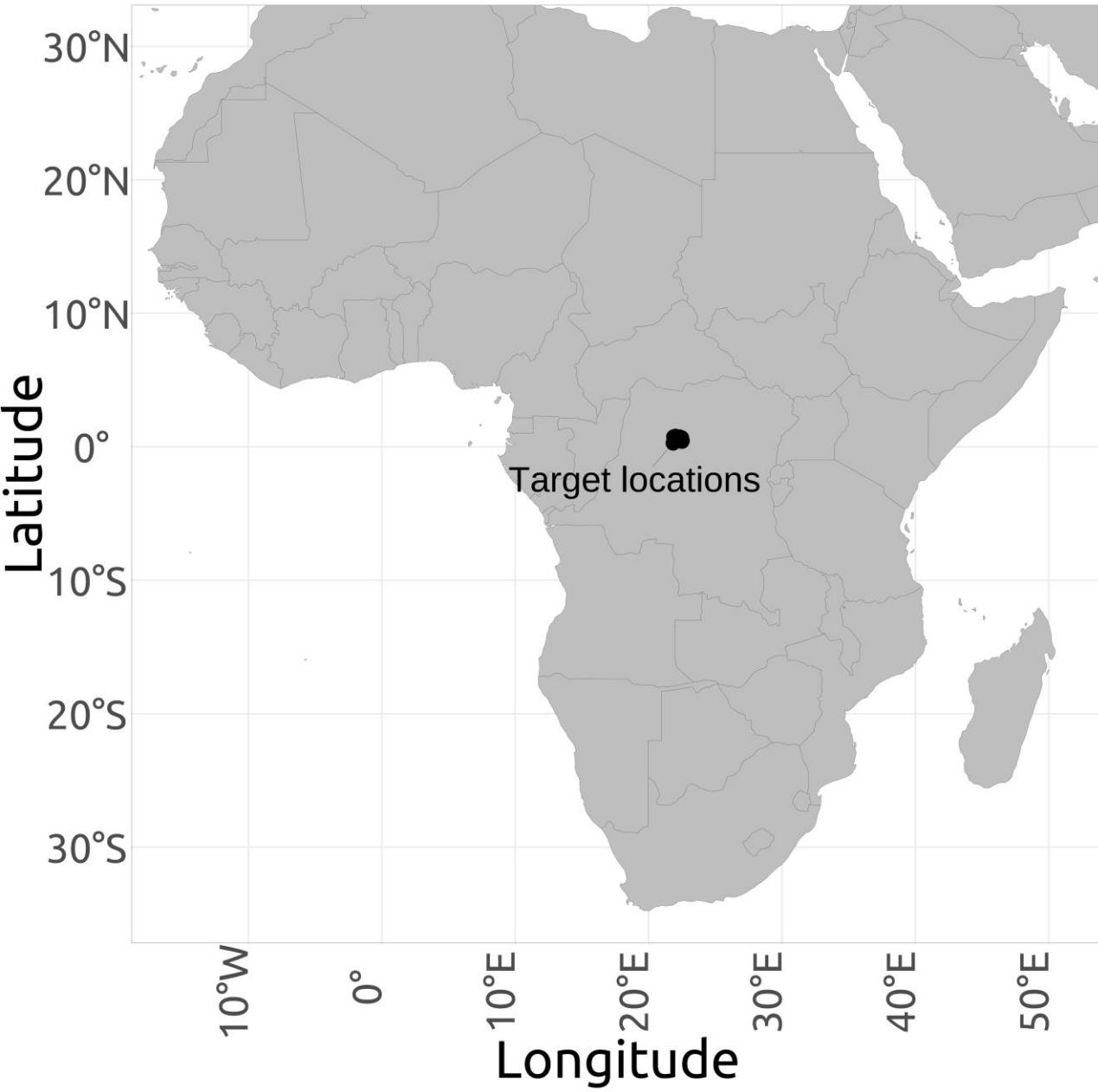


Think globally, fit locally*!

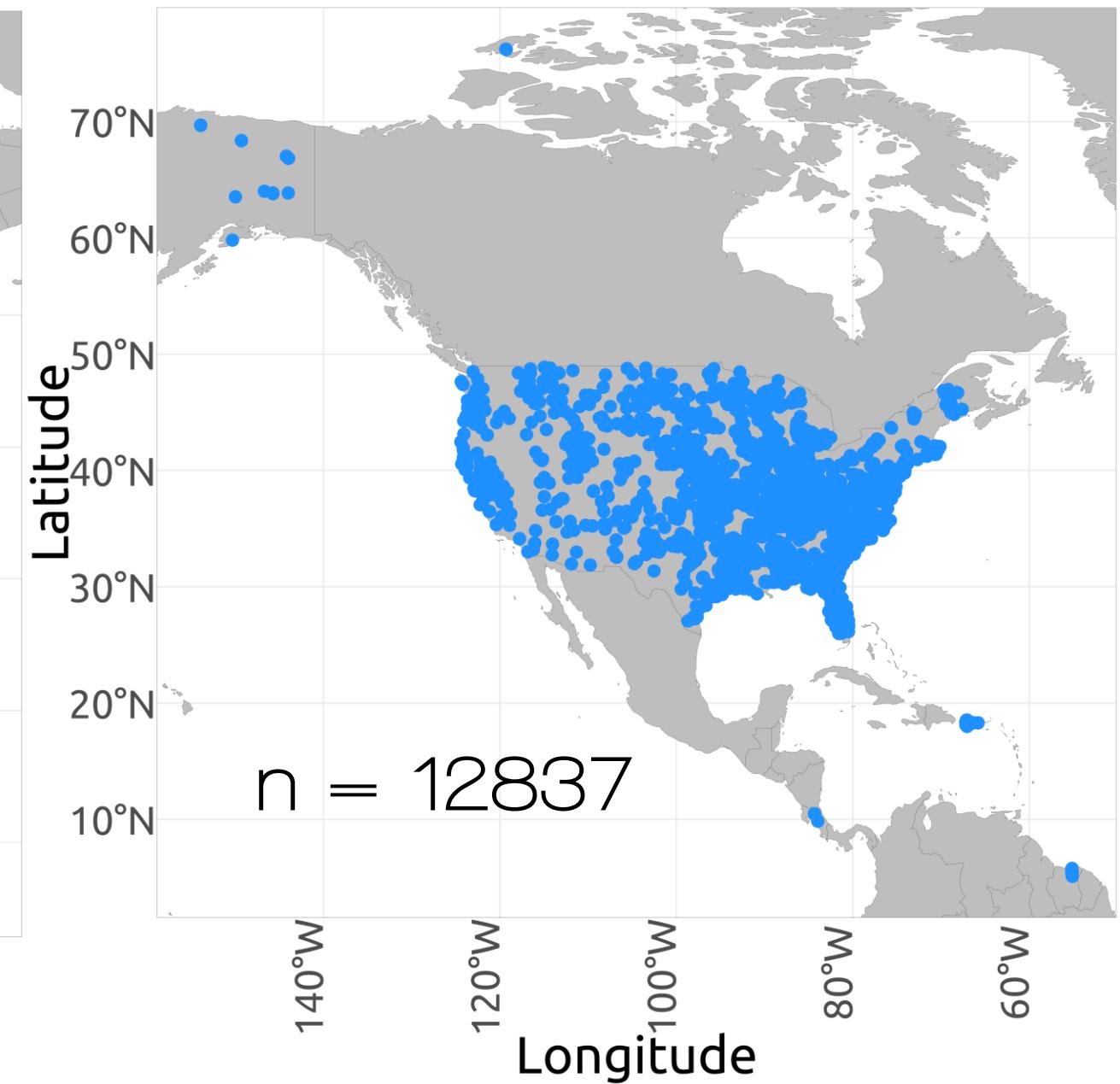
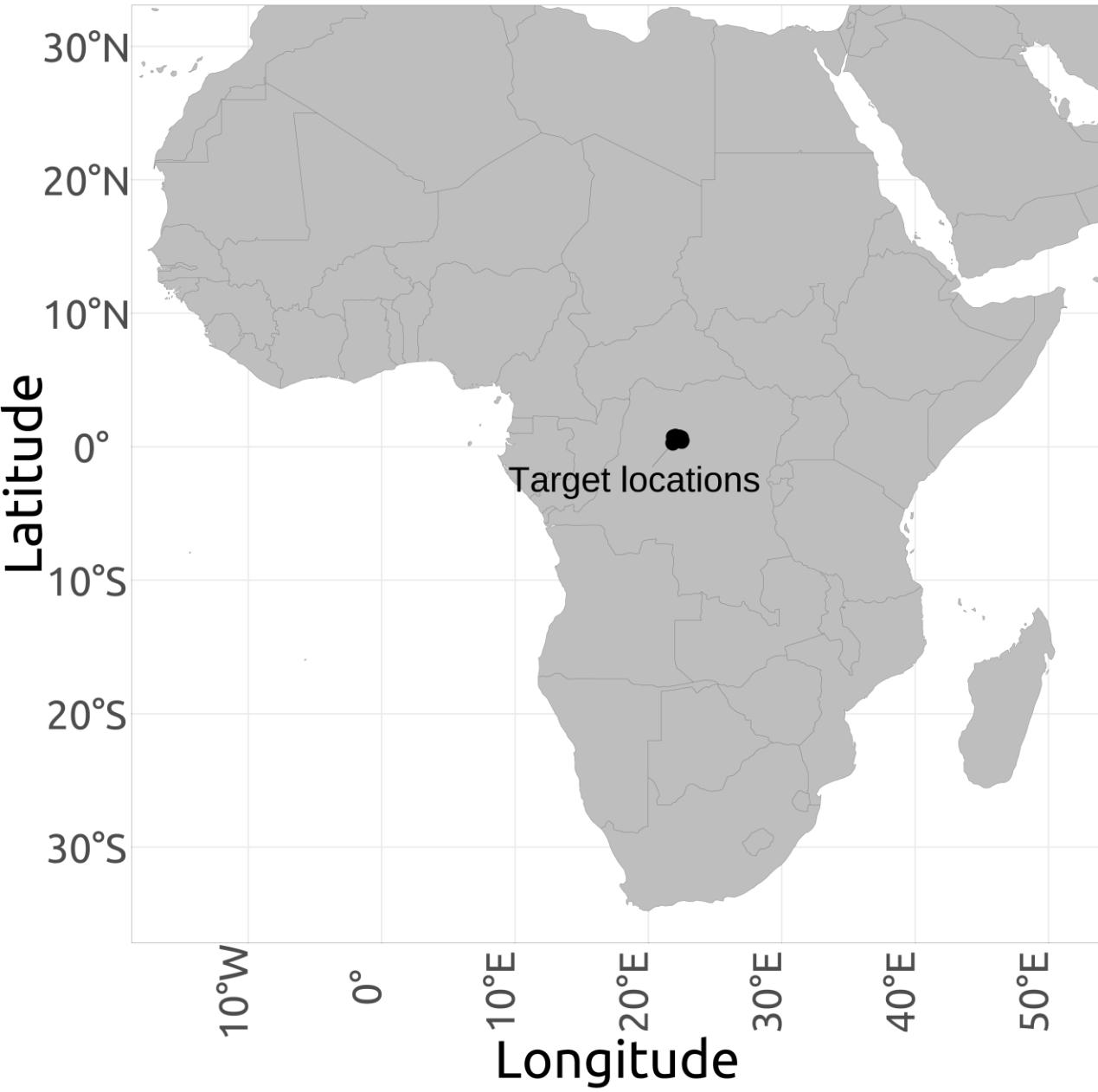
*Saul, L. K., & Roweis, S. T. (2003). Think globally, fit locally: unsupervised learning of low dimensional manifolds. *Journal of machine Learning research*, 4(Jun), 119-155.

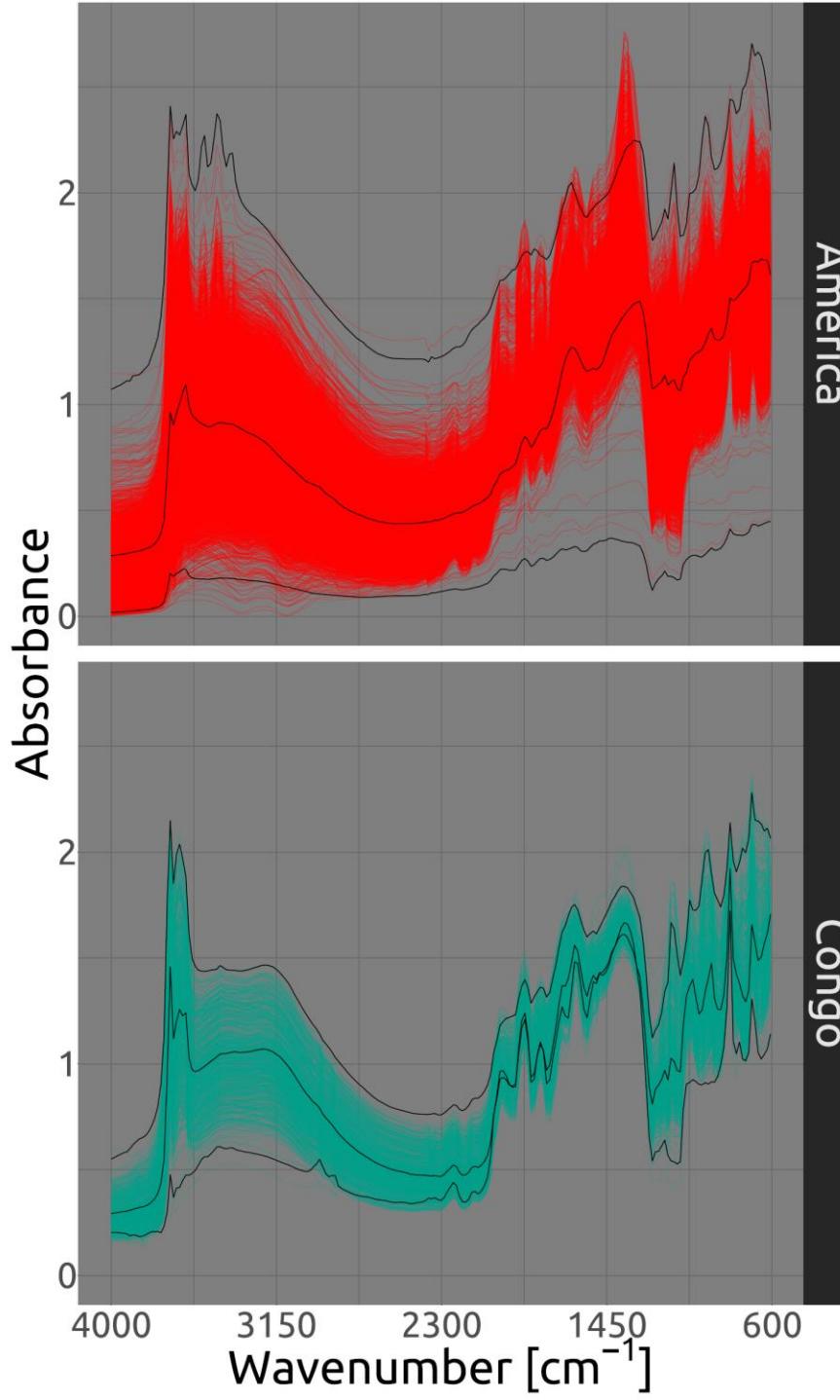
Example

Total Carbon in soils



n: 729





12837 samples

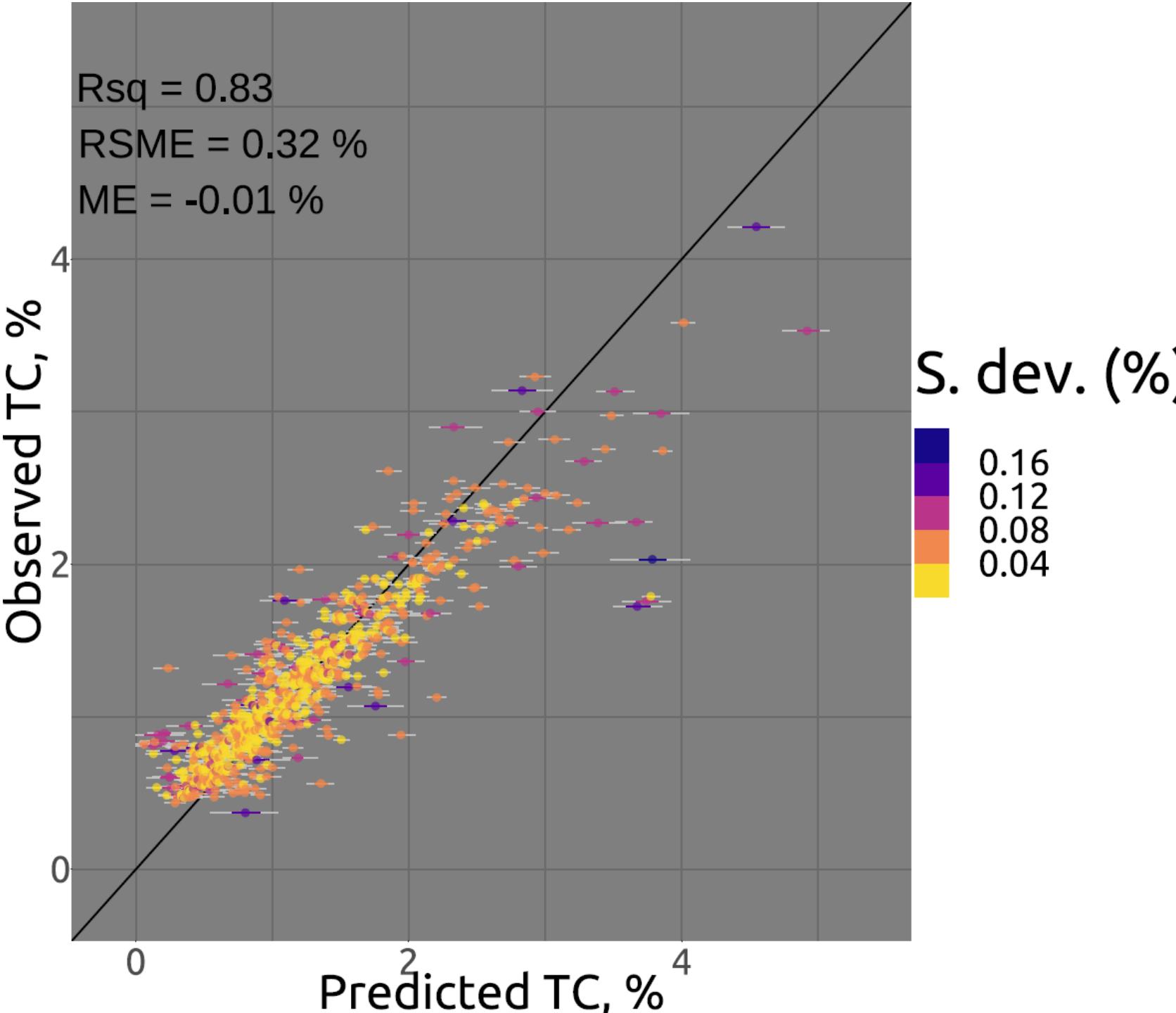
729 samples

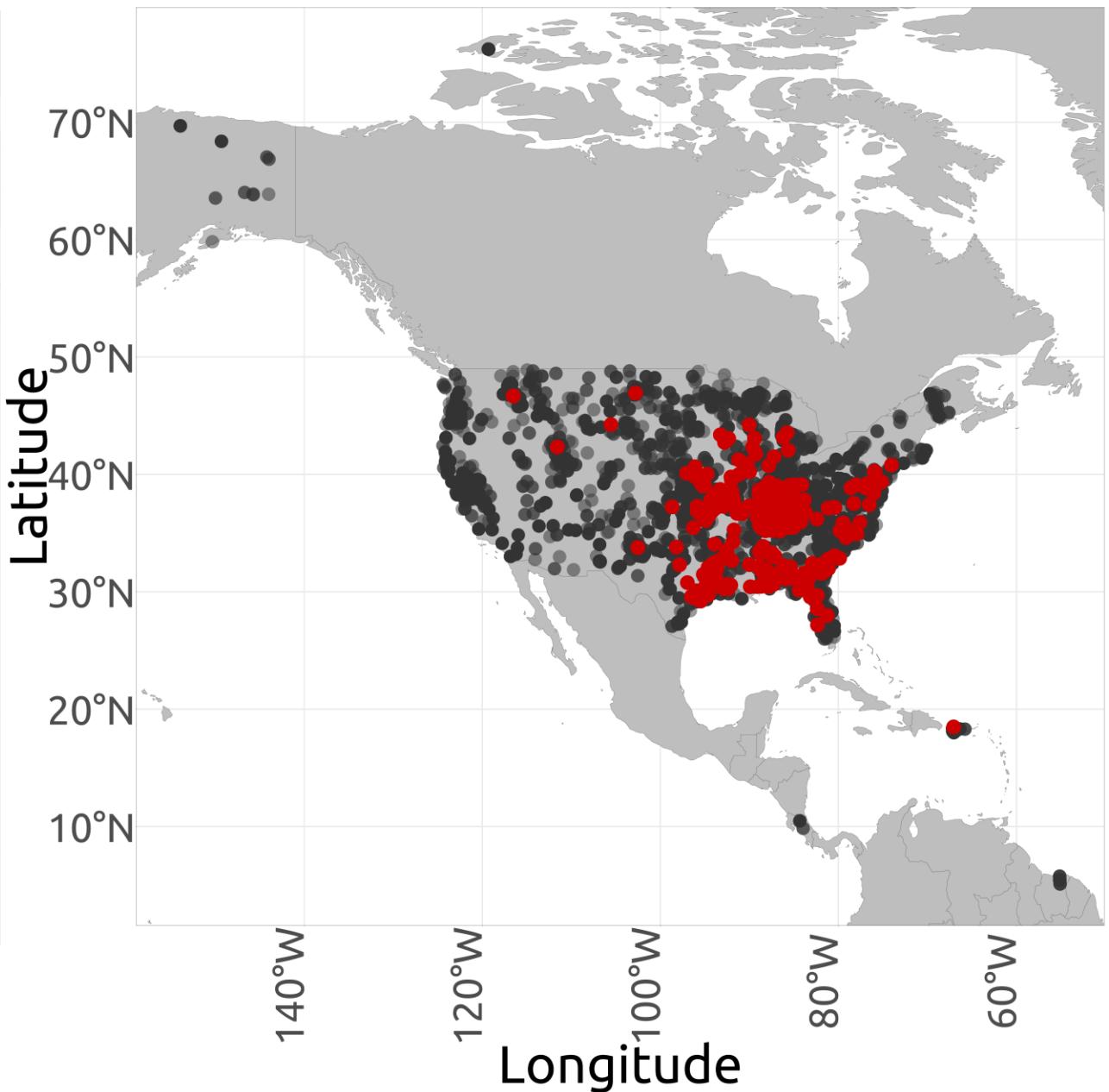
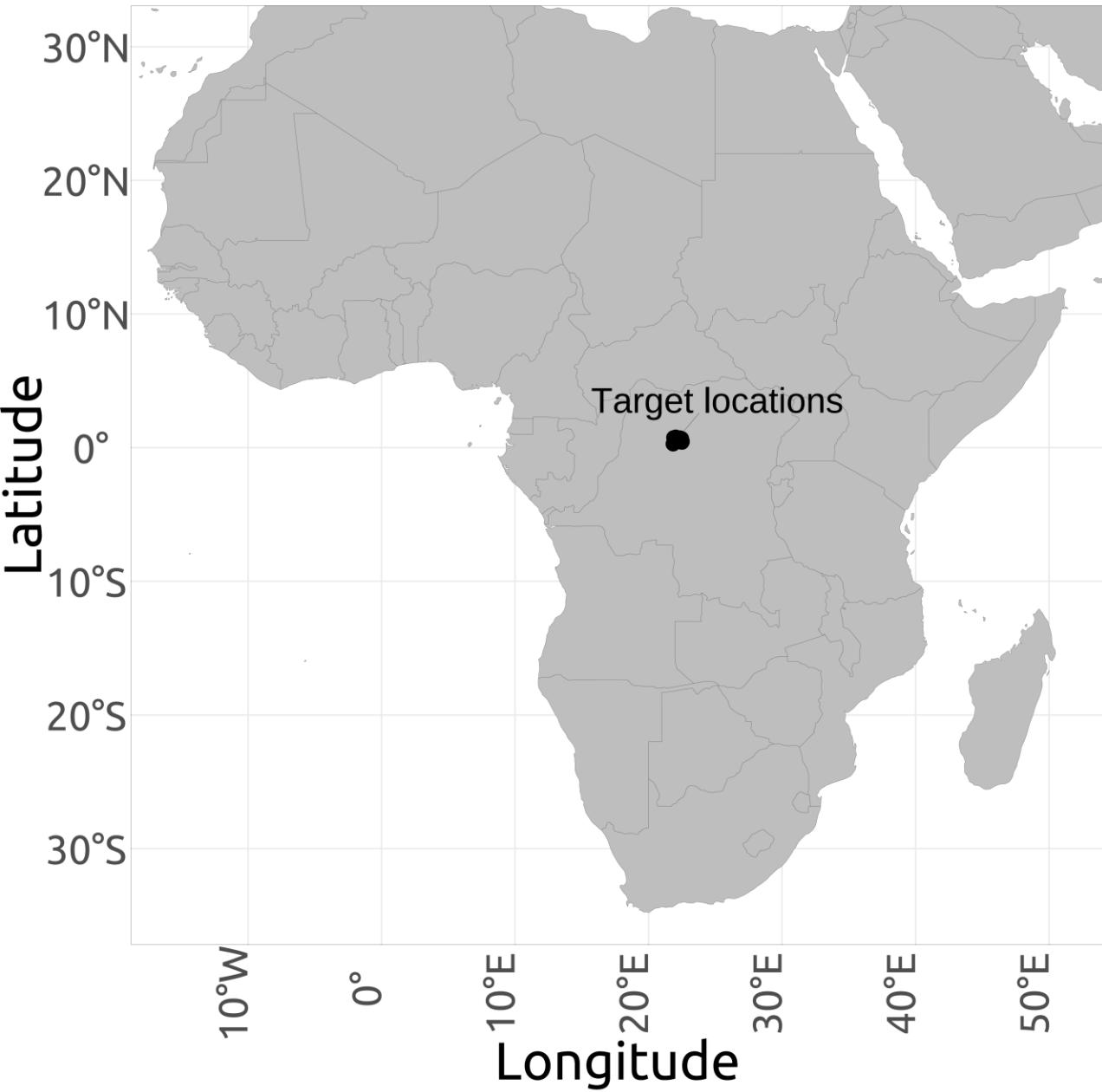
Genetically derived model

Weakness scores:

- Spectral reconstruction
- Spectral similarity
- Response accuracy 

↓
20 samples with response values

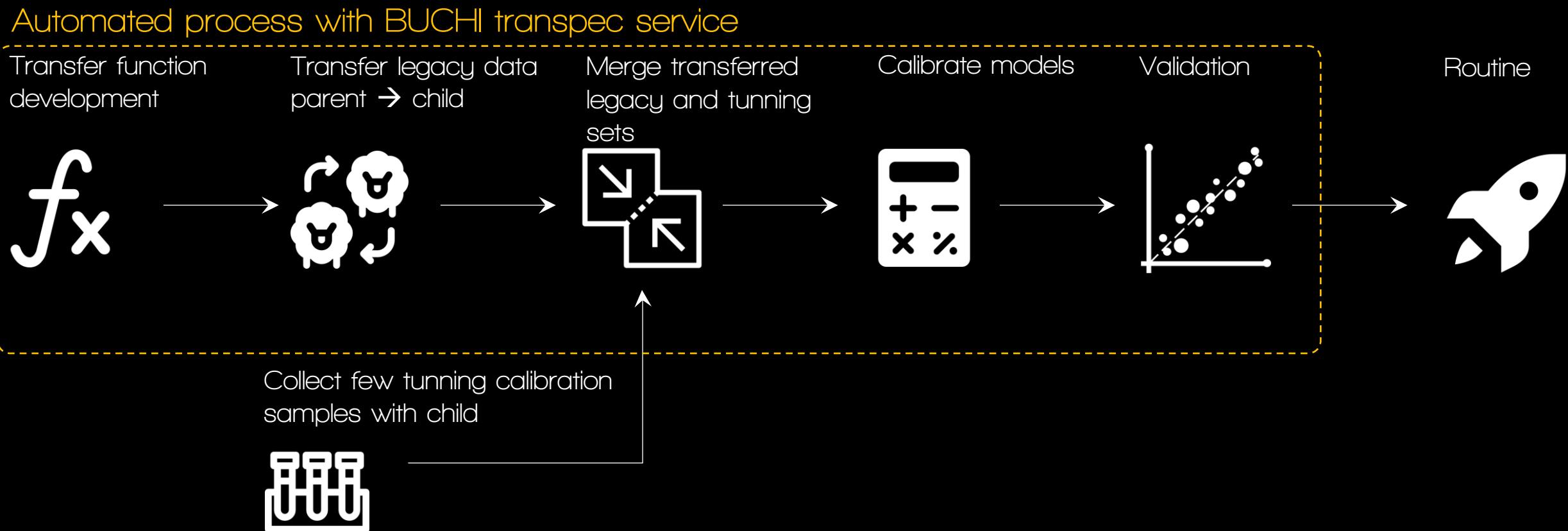




Think globally, fit locally*!

*Saul, L. K., & Roweis, S. T. (2003). Think globally, fit locally: unsupervised learning of low-dimensional manifolds. *Journal of machine learning research*, 4(Jun), 119–155.

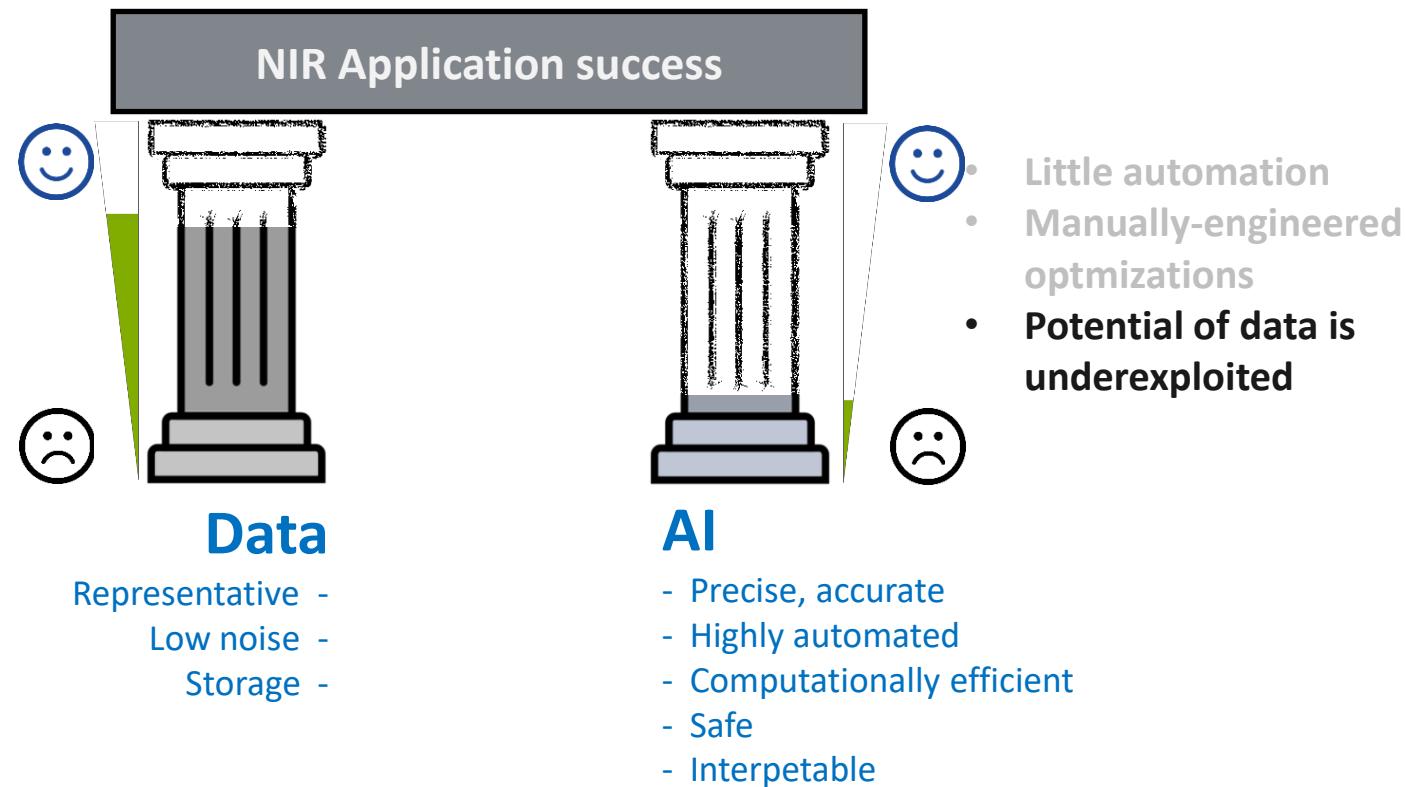
Our generalization



AI: the **what** and the **why**?

Is AI the holy grail for success?

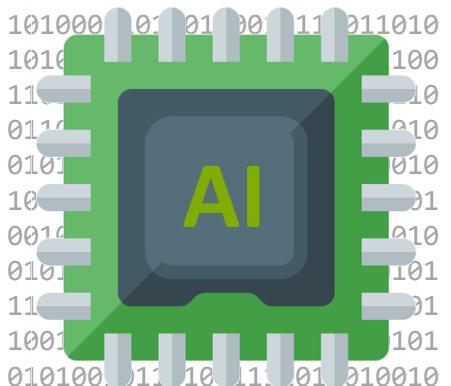
How much cement do we have to build the pillars of an excellent NIR application system?



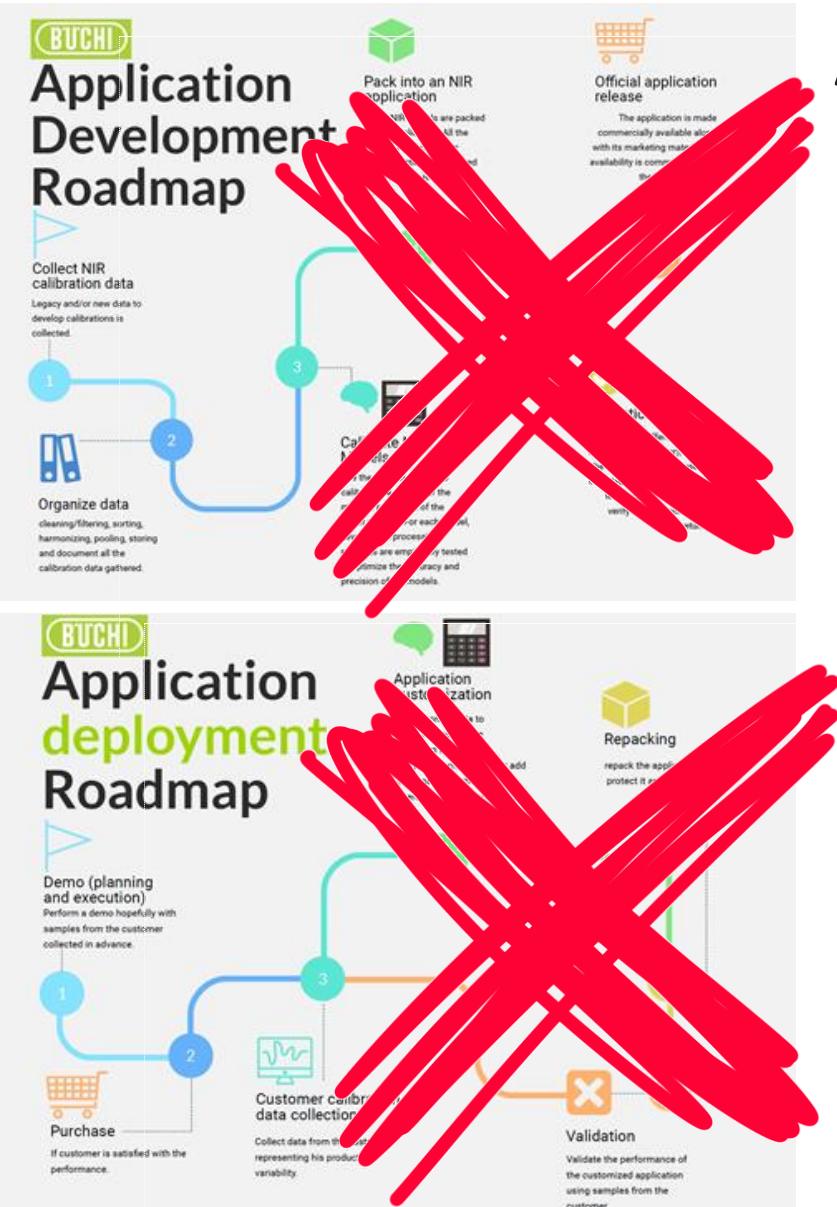
From current situation... (BUCHI NIRCal software)



To a sound NETCAL Solution...



AI: the what and the why?



We need to remove all this heavy calibration burden because:

It is expensive

AI: the what and the why?



We need to remove all this heavy calibration burden because:

- It is expensive
- Sales force should focus on sales (not on learning and dealing with chemometrics)

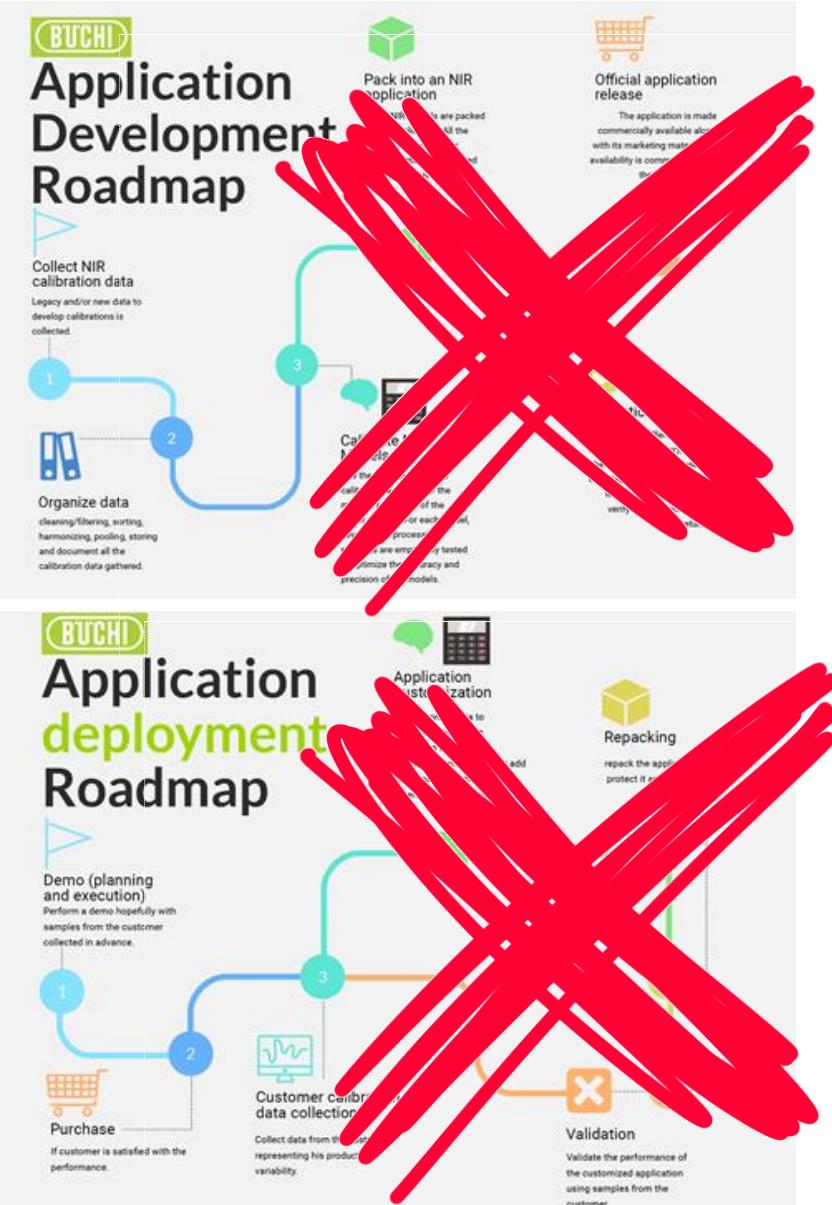
AI: the what and the why?



We need to remove all this heavy calibration burden because:

- **It is expensive**
- **Sales force should focus on sales** (not on learning and dealing with chemometrics)
- **For customers**, many time-consuming iterations to get their applications ready **create frustration.**

AI: the what and the **why**?



We need to remove all this heavy calibration burden because:

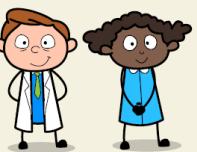
- **It is expensive**
- **Sales force should focus on sales** (not on learning and dealing with chemometrics)
- **For customers**, many time-consuming iterations to get their applications ready **create frustration**.
- **Chemometrics** is a complex subject which **requires time to be learned/taught**



Food and Agriculture
Organization of the
United Nations

GLOSOLAN is a Global
Soil Laboratory Network
which aims to harmonize soil
analysis methods and data so
that soil information is
comparable and interpretable
across laboratories, countries
and regions.

Good practices on purchasing and operating laboratory equipment



This document aims to provide laboratory staff and managers with guidance on **what to do** and **what not to do** when purchasing laboratory equipment or receiving it as a donation.

The document is divided into three sections to provide users with as much support as possible. It also includes guidance on good practices in managing consumables and hazardous substances.



For service labs...

Proposed Machine Learning Framework

1 Sampling

Customer measures few samples which are representative of its product variability. Reference property values are not mandatory. Embedded BUCHI ML-based sampling tools guide this process.



2

Upload data

User uploads all the data collected from their product.

3



Machine Learning (ML) API

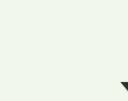
BUCHI ML algorithms are used to build optimal predictive models with selected BUCHI legacy training data that optimally match the user samples. If user data contains reference/response values, these samples are included in the models.



4

Download Application

Download the application developed by ML.



5

Validate

User checks if the performance is reaching his expectations. If not, go to step 1.





Contributions

- Understanding of the standardization concepts developed by competitors
- Build tools an internal standardization service
- Better understanding of the uncertainty sources and their impact on application performance
- Development of our own standardization concept
- Close cooperation with NIR assembly
- Development of standard operating procedures
- Implementation of incoming inspection of white references