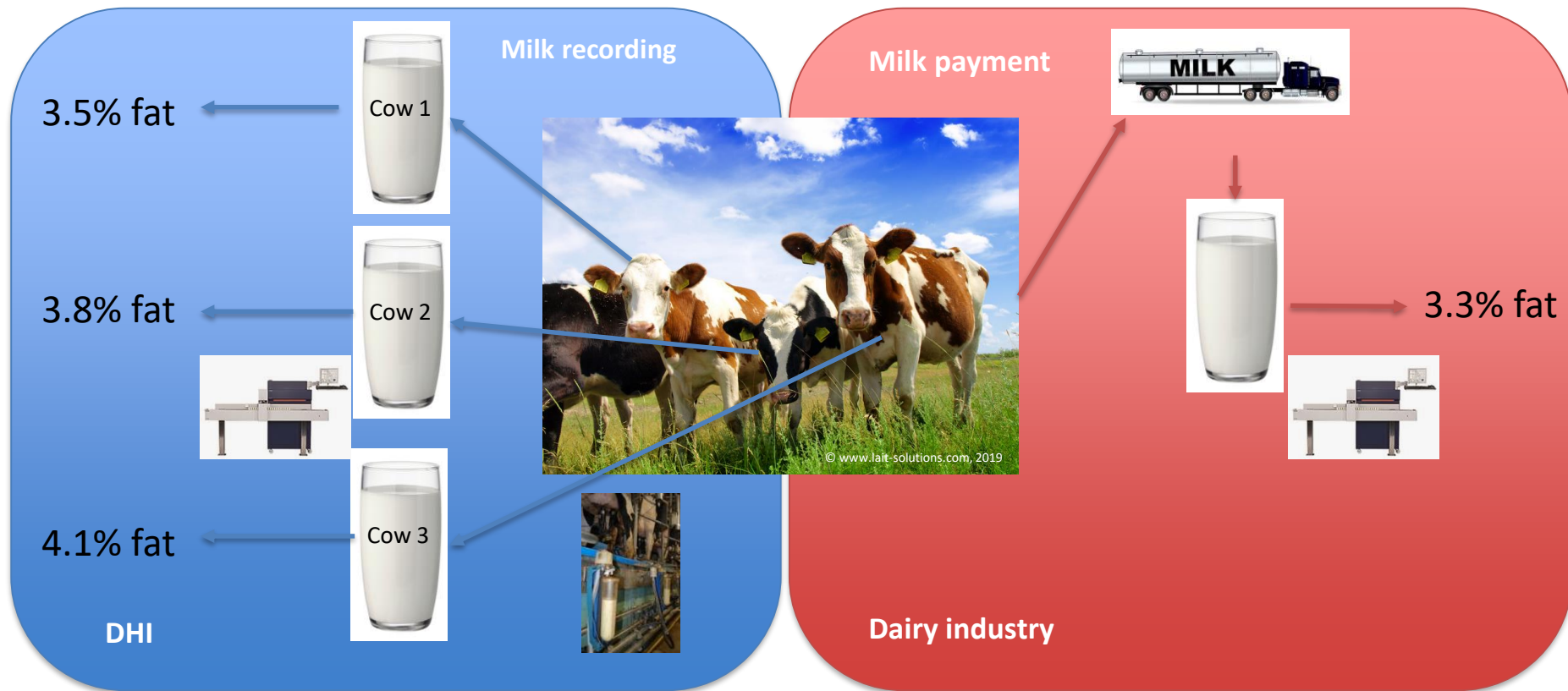LIÈGE université
**Gembloux**
**Agro-Bio Tech**

**Lei ZHANG**

Ma Y.; DEHARENG F.; Grelet C.; COLINET F.; GENGLER N.; SOYEURT H.

# Can the calculation of a spectral Global H distance ensure the quality of international based MIR predictions?

ICAR, Prague, 17/ 06/ 2019-21/06/2019
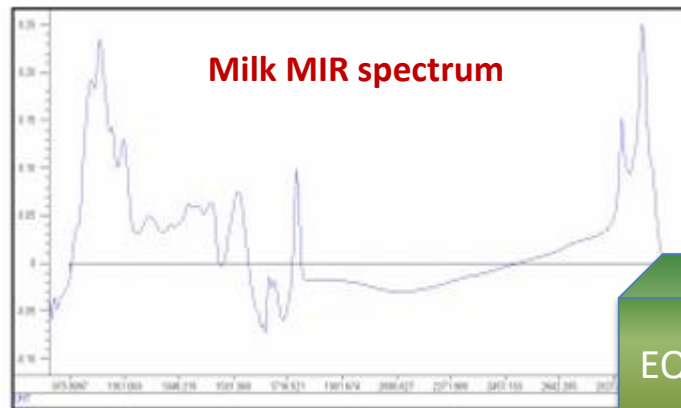
# Milk Recording Scheme

# Milk Analysis



Milk recording
(About 1 month for each cow)

%fat

**Milk MIR spectrum**
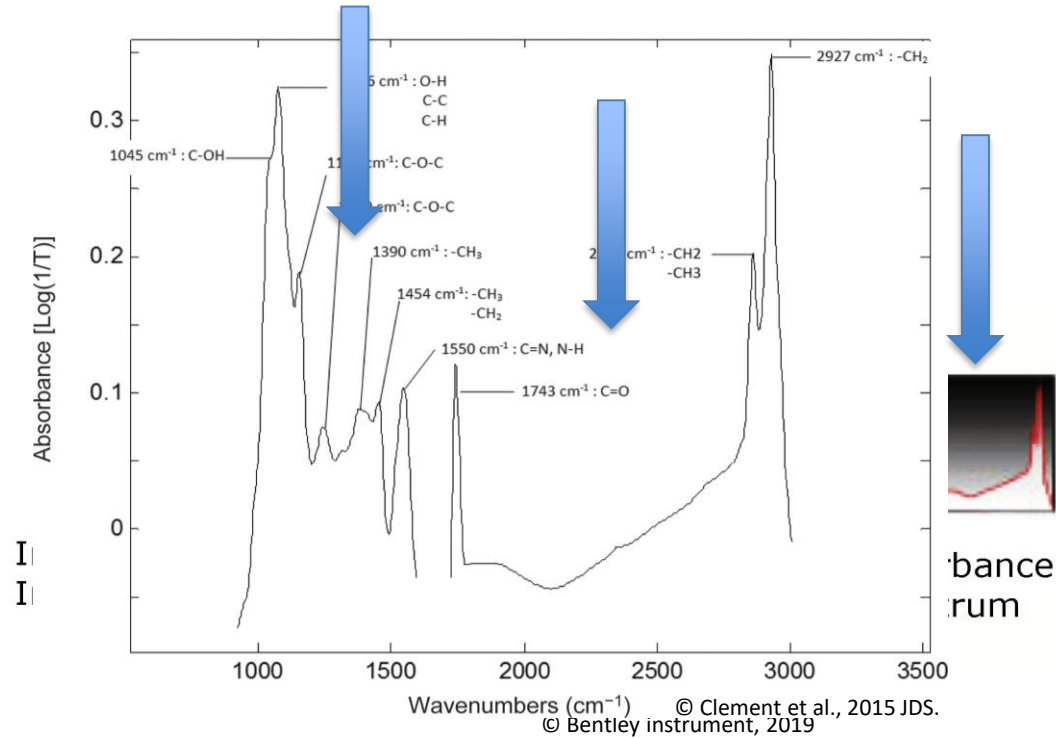
EQUATION

# What is Mid-infrared spectrum?

▶ Approximately 2,500-25,000nm (4,000-400 cm$^{-1}$)



© FOSS,2019

# Principle of MIR spectrometry



© Clement et al., 2015 JDS.
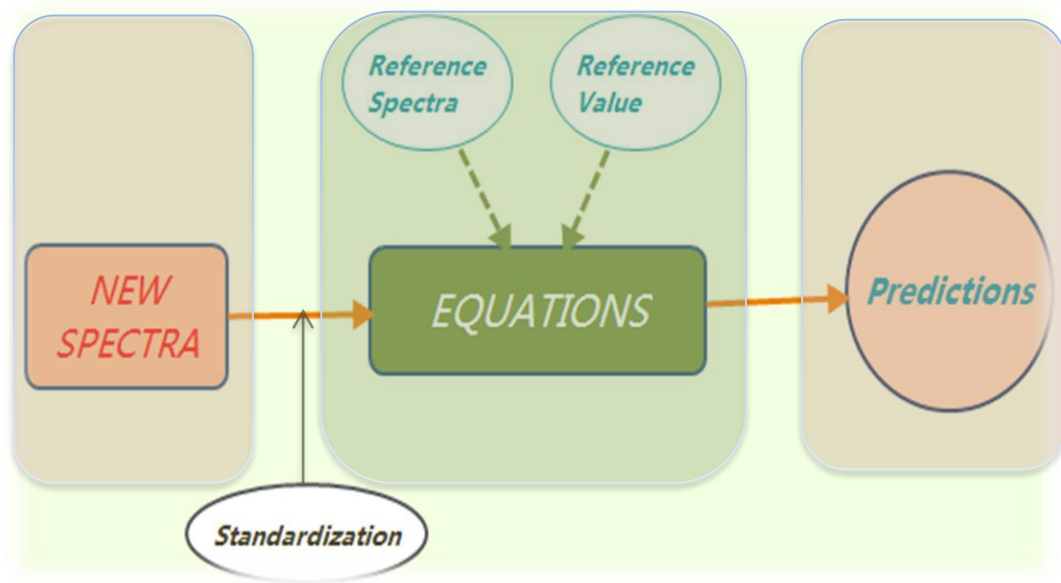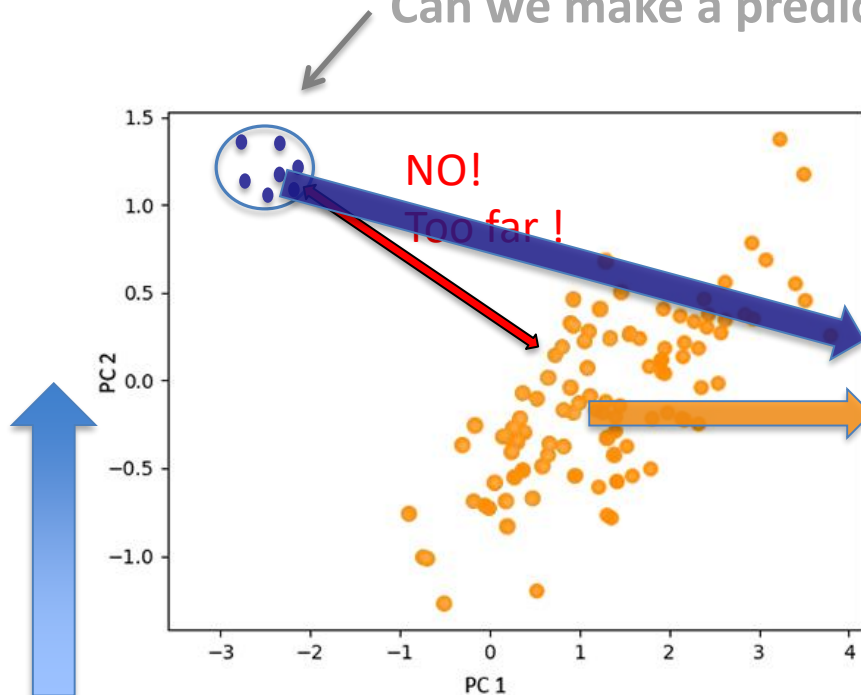© Bentley instrument, 2019

# How can we make a prediction ?

# Can we make a prediction for all spectra ?

Can we make a prediction from those spectra?



NO!
Too far !

Milk MIR Spectrum
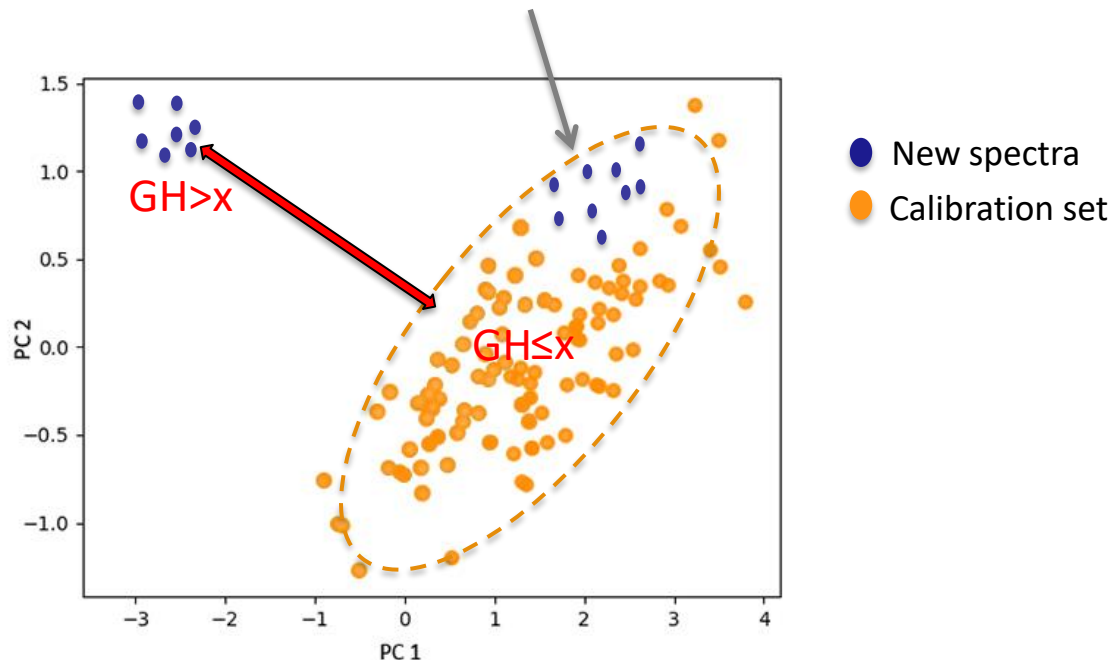
© Modified from scikit-learn. 2019

# Can we make a prediction for all spectra ?
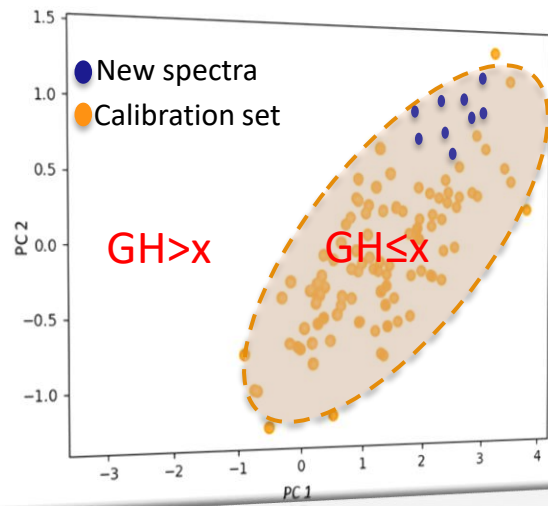
Can I make a prediction from those spectra?



GH>x

GH≤x

● New spectra
● Calibration set

Yeah

© Modified from scikit-learn. 2019

**Mahalanobis Distance:**

$$D_M(\vec{x}) = \sqrt{(\vec{x} - \vec{\mu})^T S^{-1} (\vec{x} - \vec{\mu})}$$

Where: $\vec{x}$ is PC scores of one spectrum;
$\vec{\mu}$ is the mean of PC scores of spectra in the calibration set;
$S$ is covariance matrix between PC scores of the calibration spectra
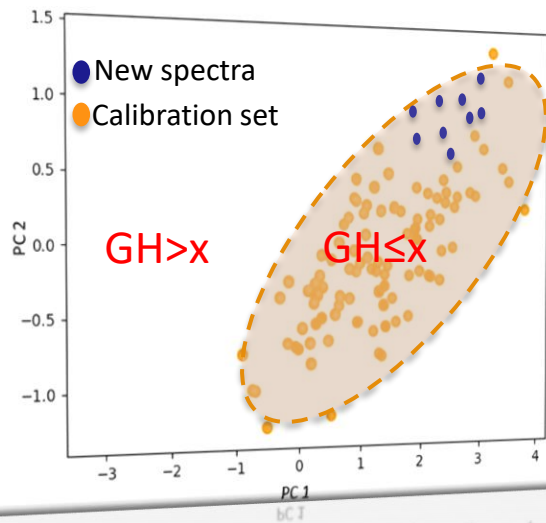
© Wikipedia
P.C. Mahalaobis
(1893 – 1972)

New spectra
Calibration set

GH>x      GH≤x

© Modified from scikit-learn. 2019

# Mahalanobis Distance:

$$D_M(\vec{x}) = \sqrt{(\vec{x} - \vec{\mu})^T S^{-1}(\vec{x} - \vec{\mu})}$$

**GH:** Global H which is the Standardized Mahalanobis Distance

$$GH = \frac{DM}{nPCs}$$

Where:
**DM** is the distance calculated from the fomular;
**nPCs** is the number of the principal components from PCA



- New spectra
- Calibration set

GH>x    GH≤x

LIÈGE université
Gembloux
Agro-Bio Tech

# What is the accuracy of prediction of international spectrum?

**Milk Recording**

MIR

STD

**MIR**$_{STD}$

**MIR**$_{Cal}$

R

Prediction

## Standardization of milk mid-infrared spectra from a European dairy network

C. Grelet[1], J.A. Fernández Pierna[1], P. Dardenne, V. Baeten, F. Dehareng[2]

Walloon Agricultural Research Center, Valorisation of Agricultural Products Department, 24 Chaussée de Namur, 5030 Gembloux, Belgium
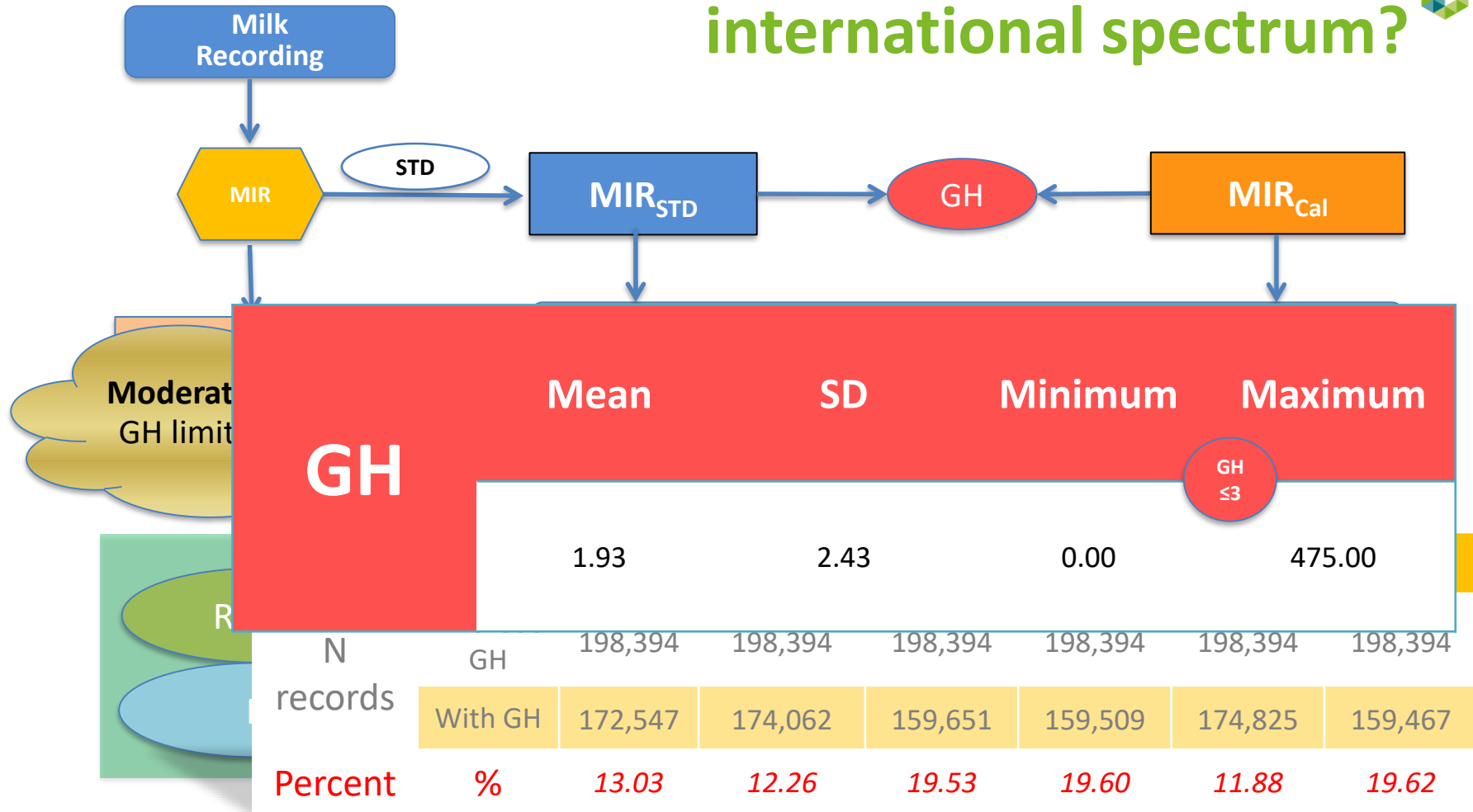
# What is the accuracy of prediction of international spectrum?

Milk Recording

MIR

STD

MIR$_{STD}$

GH

MIR$_{Cal}$

**Moderate** GH limit

GH ≤3

| GH | Mean | SD | Minimum | Maximum | | |
|---|---|---|---|---|---|---|
| | 1.93 | 2.43 | 0.00 | 475.00 | | |
| N records | GH | 198,394 | 198,394 | 198,394 | 198,394 | 198,394 | 198,394 |
| | With GH | 172,547 | 174,062 | 159,651 | 159,509 | 174,825 | 159,467 |
| Percent | % | *13.03* | *12.26* | *19.53* | *19.60* | *11.88* | *19.62* |

# Descriptive statistics

Table 1. Descriptive statistics of predicted value

| Traits g/dL | Reference value | | Predicted value | | Predicted value (GH <= 3) | |
|---|---|---|---|---|---|---|
| | Mean | SD | Mean | SD | Mean | SD |
| Fat | 3.97 | 0.95 | 3.99 | 0.95 | 3.90 | 0.86 |
| Protein | 3.43 | 0.40 | 3.53 | 0.46 | 3.52 | 0.39 |
| MFA | 0.86 | 0.27 | 1.15 | 0.36 | | 0.31 |
| PFA | 0.07 | 0.04 | 0.15 | 0.05 | 0.15 | 0.04 |
| SFA | 2.62 | 0.67 | 2.64 | 0.68 | 2.59 | 0.62 |
| UFA | 0.93 | 0.31 | 1.29 | 0.39 | 1.25 | 0.34 |

# The correlation coefficient



Why?

# Squared residual and GH

# GH limitation decreased RMSE for most traits

# Conclusion:

▶ GH limitation helps to **ensure the quality** of the MIR predictions

▶ It allows **avoiding spectral extrapolation**

▶ More work needed to be done to get **more accurate** predictions...

**Predicted Accuracy**

**Extrapolation Control**

GH limitation

Ma Y.; DEHARENG F.; Grelet C.; COLINET F.; GENGLER N.; SOYEURT H.,&

**Lei ZHANG**

**Email:** lei.zhang@doct.uliege.be

ICAR, Prague, 17/ 06/ 2019-21/06/2019

# Additional information:
## Why do PCA?

▶ To decrease the dimensionality of the raw data

▶ To make it easy for calculating the inverse of the covariance matrix



Lever et al., 2017 Nature Method 2017

# Additional information:

## Why GH ≤ 3?



68-95-99.7 Rule