



HAL
open science

Model and auxiliary data for an accurate estimate of the field mean: Appendix

Baptiste Oger, Sébastien Roux, Gilles Le Moguedec, Bruno Tisseyre

► **To cite this version:**

Baptiste Oger, Sébastien Roux, Gilles Le Moguedec, Bruno Tisseyre. Model and auxiliary data for an accurate estimate of the field mean: Appendix. 2021. <hal-03099451>

HAL Id: hal-03099451

<https://hal.science/hal-03099451v1>

Submitted on 27 Feb 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

Model and auxiliary data for an accurate estimate of the field mean: Appendix

B. Oger^{1,2}, S. Roux², G. Le Moguedec³ and B. Tisseyre¹

¹ITAP, Univ. Montpellier, Montpellier SupAgro, INRAE, France

²MISTEA, Univ. Montpellier, Montpellier SupAgro, INRAE, France

³AMAP, Univ. Montpellier, INRAE, France

baptiste.oger@supagro.fr

This document supplements the proposed paper ‘*Model and auxiliary data for an accurate estimate of the field mean*’ for the 13th European Conference on Precision Agriculture - ECPA 2021 (Oger et al. 2021). This conference paper proposes a discussion on inference methods for crop production estimation following previous studies (Carillo et al., 2016; Araya-Alma et al. 2019). This appendix includes the full demonstration of the theoretical results. For context and interpretation, please refer to the conference paper.

Notations and hypothesis

For a given plot, K is defined as the within-plot set of sites, numbered form 1 to k . For a given site $i \in K$, X_i and Y_i are respectively the auxiliary variable and the variable of interest measured on that site. An auxiliary data is assumed available at low cost and at a high spatial resolution. The auxiliary variable X_i is known for every site but the variable of interest Y_i only for a set N of sampled sites (with $Cardinal(N) = n$). It is assumed that Y_i can be predicted everywhere form X_i via a simple linear model:

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i \quad \text{with} \quad \varepsilon_i \sim N(0, \sigma^2) \quad (1)$$

Let \bar{Y} be the average value of Y for the whole plot:

$$\bar{Y} = \frac{1}{k} \times \sum_{i \in K} Y_i \quad (2)$$

To focus on estimator related error, two hypotheses are made here. First, the calculations are made under the assumption that measurement sites are determined (the N set is fixed). Second, model residuals are considered independent, i.e. unaffected by the spatial structure of the plot.

Mean-based estimator

Estimator

A first estimator \widehat{Y}_1 of \bar{Y} based on the mean inference. \widehat{Y}_1 is defined in Equation 3 as the average of the values Y_i taken on the n sampled sites ($i \in N$):

$$\widehat{Y}_1 = \frac{1}{n} \times \sum_{i \in N} Y_i \quad (3)$$

Expectation

Main statistical properties of this estimator are given by its expected value (for the bias) and the variance (for the precision). The \widehat{Y}_1 expectation (E) is written:

$$E[\widehat{Y}_1] = E\left[\frac{1}{n} \times \sum_{i \in N} Y_i\right]$$

Although the mean estimator is not based on auxiliary data, it is assumed that such data exist to characterize its properties. By taking the linear relationship between Y_i and X_i presented in Equation 1, it is possible to replace Y_i by its expression from X_i :

$$\begin{aligned} E[\widehat{Y}_1] &= E\left[\frac{1}{n} \times \sum_{i \in N} (\beta_0 + \beta_1 X_i + \varepsilon_i)\right] \\ E[\widehat{Y}_1] &= E\left[\frac{1}{n} \times \sum_{i \in N} \beta_0\right] + E\left[\frac{1}{n} \times \sum_{i \in N} (\beta_1 X_i)\right] + E\left[\frac{1}{n} \times \sum_{i \in N} \varepsilon_i\right] \\ E[\widehat{Y}_1] &= \beta_0 + \beta_1 \overline{X_{i \in N}} + 0 \\ E[\widehat{Y}_1] &= \beta_0 + \beta_1 (\overline{X_{i \in N}} + \bar{X} - \bar{X}) \\ E[\widehat{Y}_1] &= \beta_0 + \beta_1 \bar{X} + \beta_1 (\overline{X_{i \in N}} - \bar{X}) \end{aligned}$$

And finally:

$$E[\widehat{Y}_1] = \bar{Y} + \beta_1 (\overline{X_N} - \bar{X}) \quad (4)$$

\widehat{Y}_1 presents a bias which depend on β_1 , the slope parameter, $\overline{X_N}$, the auxiliary data mean value for the sample and \bar{X} , the auxiliary data mean value for the whole plot. The difference between $\overline{X_N}$ and \bar{X} reflects the influence of sample representativeness, ideally the sample should have properties close to those of the population. Since representativeness is expressed in terms of auxiliary data, its influence depends on β_1 which characterizes the strength of the correlation between the variables Y_i and X_i .

Variance:

Looking at its variance:

$$Var[\widehat{Y}_1] = Var\left[\frac{1}{n} \times \sum_{i \in N} Y_i\right]$$

Using Equation 1:

$$Var[\widehat{Y}_1] = Var\left[\frac{1}{n} \times \sum_{i \in N} (\beta_0 + \beta_1 X_i + \varepsilon_i)\right]$$

As β_0 , β_1 and X_i correspond to numerical values, their variances are equal to 0. For X_i , this is true under the hypothesis that measurement sites are determined.

$$Var[\widehat{Y}_1] = Var\left[\frac{1}{n} \times \sum_{i \in N} \varepsilon_i\right]$$

With $Var(aX) = a^2 Var(X)$:

$$Var[\widehat{Y}_1] = \frac{1}{n^2} \times Var\left[\sum_{i \in N} \varepsilon_i\right]$$

And finally:

$$\begin{aligned} \text{Var}[\widehat{Y}_1] &= \frac{1}{n^2} \times n\sigma^2 \\ \text{Var}[\widehat{Y}_1] &= \frac{\sigma^2}{n} \end{aligned} \quad (5)$$

The variance can be described with σ , the residual variance of the linear model and n , the number of sampling sites. This result is logical, the higher the variability of the variable of interest, the more variable the estimate will be. Conversely, a larger sample size improves confidence in the estimate.

Model-based estimator

Estimator

Let \widehat{Y}_i be the estimate for Y_i based on the linear model calibrated using the values of Y_i on the n sampling sites. A second estimator, \widehat{Y}_2 , can be described from the set of values predicted by the model:

$$\widehat{Y}_2 = \frac{1}{k} \times \sum_{i=1}^k \widehat{Y}_i \quad (6)$$

which can be rewritten:

$$\begin{aligned} \widehat{Y}_2 &= \frac{1}{k} \times \sum_{i=1}^k (\widehat{\beta}_0 + \widehat{\beta}_1 X_i) \\ \widehat{Y}_2 &= \widehat{\beta}_0 + \widehat{\beta}_1 \bar{X} \end{aligned}$$

$\widehat{\beta}_0$ and $\widehat{\beta}_1$ are the respective estimators of β_0 and β_1 . In the simple linear model, these estimators are well known (Wasserman, 2004):

$$\begin{aligned} \widehat{\beta}_1 &= \frac{\sum_{i \in N} (X_i - \bar{X}_N)(Y_i - \bar{Y}_N)}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \\ \widehat{\beta}_0 &= \bar{Y}_N - \widehat{\beta}_1 \times \bar{X}_N \end{aligned}$$

And where \bar{X}_N represents the mean of X_i for $i \in N$ and \bar{Y}_N represents the mean Y_i for $i \in N$.

Expectation

As for \widehat{Y}_1 , the expectation expression is derived for \widehat{Y}_2 :

$$\begin{aligned} E[\widehat{Y}_2] &= E[\widehat{\beta}_0 + \widehat{\beta}_1 \bar{X}] \\ E[\widehat{Y}_2] &= E[\widehat{\beta}_0] + E[\widehat{\beta}_1] \bar{X} \end{aligned}$$

One of the properties of the linear model is that the estimators $\widehat{\beta}_0$ and $\widehat{\beta}_1$ for β_0 and β_1 are unbiased (Wasserman, 2004). In other words, that $E[\widehat{\beta}_0] = \beta_0$ and $E[\widehat{\beta}_1] = \beta_1$. As a result:

$$E[\widehat{Y}_2] = \beta_0 + \beta_1 \bar{X}$$

And finally:

$$\mathbf{E}[\widehat{Y}_2] = \bar{Y} \quad (7)$$

This result shows that the estimation based on model and auxiliary data inference is unbiased.

Variance:

To be able to fully compare the two estimators, it is nevertheless necessary to look at the \widehat{Y}_2 variance:

$$Var[\widehat{Y}_2] = Var[\widehat{\beta}_0 + \widehat{\beta}_1 \bar{X}]$$

The variance of a sum can be decomposed as follows: $Var(X + Y) = Var(X) + Var(Y) + 2Cov(X, Y)$. As a result:

$$Var[\widehat{Y}_2] = Var[\widehat{\beta}_0] + Var[\widehat{\beta}_1 \bar{X}] + 2Cov[\widehat{\beta}_0, \widehat{\beta}_1 \bar{X}]$$

$$Var[\widehat{Y}_2] = \sigma_{\widehat{\beta}_0}^2 + \sigma_{\widehat{\beta}_1}^2 \bar{X}^2 + 2\bar{X}Cov[\widehat{\beta}_0, \widehat{\beta}_1]$$

However, in the simple linear model, the expressions of the components $\sigma_{\widehat{\beta}_0}^2$, $\sigma_{\widehat{\beta}_1}^2$ and $Cov[\widehat{\beta}_0, \widehat{\beta}_1]$ are known and expressed (Wasserman, 2004):

$$\begin{aligned} \sigma_{\widehat{\beta}_1}^2 &= \frac{\sigma^2}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \\ \sigma_{\widehat{\beta}_0}^2 &= \sigma^2 \times \left(\frac{1}{n} + \frac{\bar{X}_N^2}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \right) \\ Cov[\widehat{\beta}_0, \widehat{\beta}_1] &= -\frac{\bar{X}_N}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \times \sigma^2 \end{aligned}$$

By replacing these components by their expressions in the expression of \widehat{Y}_2 variance:

$$\begin{aligned} Var[\widehat{Y}_2] &= \sigma^2 \times \left(\frac{1}{n} + \frac{\bar{X}_N^2}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \right) + \frac{\sigma^2}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \bar{X}^2 - \frac{\bar{X}_N}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \times \sigma^2 2\bar{X} \\ Var[\widehat{Y}_2] &= \frac{\sigma^2}{n} + \frac{\bar{X}_N^2}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \sigma^2 + \frac{\sigma^2}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \bar{X}^2 - \frac{\bar{X}_N}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \times \sigma^2 2\bar{X} \end{aligned}$$

It is then factorized by $\frac{\sigma^2}{\sum_{i \in N} (X_i - \bar{X}_N)^2}$:

$$\begin{aligned} Var[\widehat{Y}_2] &= \frac{\sigma^2}{n} + \frac{\sigma^2}{\sum_{i \in N} (X_i - \bar{X}_N)^2} (\bar{X}_N^2 + \bar{X}^2 - 2\bar{X} \times \bar{X}_N) \\ \mathbf{Var}[\widehat{Y}_2] &= \frac{\sigma^2}{n} + \frac{(\bar{X}_N - \bar{X})^2}{\sum_{i \in N} (X_i - \bar{X}_N)^2} \times \sigma^2 \end{aligned} \quad (8)$$

As for \widehat{Y}_1 , the \widehat{Y}_2 variance depends on σ , the residual variance of the linear model and n , the number of sampling sites. It also depends on the sampled sites and their auxiliary data. To reduce its variance: the sample must i) be representative, so its mean (\bar{X}_N) is close to the field mean (\bar{X}) and ii) present a strong dispersion so the samples sites are far apart from their own mean $\sum_{i \in N} (X_i - \bar{X}_N)^2$. This is due to the fact that a dispersed sample will allow a better estimation of the model parameters.

Root mean squared error

By looking at Equations 4 and 7, \widehat{Y}_2 appears biased while \widehat{Y}_1 is unbiased. But at the same time, \widehat{Y}_2 has higher variance than \widehat{Y}_1 (Equations 5 and 8). To be able to compare them, theoretical Mean Square Error (MSE) and Root Mean Square Error (RMSE) are computed for both estimators from bias and variance as stated in Equation 9:

$$MSE(\widehat{Y}) = Bias(\widehat{Y})^2 + Variance(\widehat{Y}) \quad (9)$$

$$RMSE(\widehat{Y}) = \sqrt{MSE(\widehat{Y})} \quad (10)$$

Based on Equation 10, RMSE of \widehat{Y}_1 and \widehat{Y}_2 are given respectively in Equation 11 and Equation 12:

$$RMSE(\widehat{Y}_1) = \sqrt{\frac{\sigma^2}{n} + \beta_1^2 \times (\overline{X_{i \in N}} - \bar{X})^2} \quad (11)$$

$$RMSE(\widehat{Y}_2) = \sqrt{\frac{\sigma^2}{n} + \frac{(\overline{X_N} - \bar{X})^2}{\sum_{i \in N} (X_i - \overline{X_N})^2} \times \sigma^2} \quad (12)$$

Derive the expression of $RMSE(\widehat{Y}_2) < RMSE(\widehat{Y}_1)$ allows to identify which condition estimator \widehat{Y}_2 is provides lower RMSE than \widehat{Y}_1 and thus when estimation by the linear model should be preferred.

$$\begin{aligned} RMSE(\widehat{Y}_2) < RMSE(\widehat{Y}_1) &\Rightarrow MSE(\widehat{Y}_2) < MSE(\widehat{Y}_1) \\ MSE(\widehat{Y}_2) < MSE(\widehat{Y}_1) &\Rightarrow \left[\frac{\sigma^2}{n} + \frac{(\overline{X_N} - \bar{X})^2}{\sum_{i \in N} (X_i - \overline{X_N})^2} \times \sigma^2 \right] < \left[\frac{\sigma^2}{n} + \beta_1^2 \times (\overline{X_{i \in N}} - \bar{X})^2 \right] \\ &\Rightarrow \frac{(\overline{X_N} - \bar{X})^2}{\sum_{i \in N} (X_i - \overline{X_N})^2} \times \sigma^2 - \beta_1^2 \times (\overline{X_{i \in N}} - \bar{X})^2 < 0 \\ &\Rightarrow (\overline{X_N} - \bar{X})^2 \times \left[\frac{\sigma^2}{\sum_{i \in N} (X_i - \overline{X_N})^2} - \beta_1^2 \right] < 0 \end{aligned}$$

From these expressions, $RMSE(\widehat{Y}_2) < RMSE(\widehat{Y}_1)$ if and only if :

$$\frac{\sigma^2}{\beta_1^2} < \sum_{i \in N} (X_i - \overline{X_N})^2 \quad (13)$$

Following this result, the choice of the inference strategy must be done according to i) σ^2 , the residual variance of the linear model; ii) β_1 , the regression slope; iii) $\sum_{i \in N} (X_i - \overline{X_N})^2$, the dispersion of the sample in terms of auxiliary data. The $\frac{\beta_1^2}{\sigma^2}$ ratio reflects the existing correlation between the estimated and the auxiliary data. If the regression slope is small in front of the residuals variance, the correlation is low and it favours the use of a mean for estimator. Conversely, if the slope is strong in front of σ and if the sample values are dispersed (for a better estimate of the model parameters), this favours the use of an estimator based on the linear model.

References:

- Araya-Alman, M., Leroux, C., Acevedo-Opazo, C., Guillaume, S., Valdés-Gómez, H., Verdugo-Vásquez, N., Pañitrur-De la Fuente, C., Tisseyre, B. (2019). A new localized sampling method to improve grape yield estimation of the current season using yield historical data. *Precision Agriculture*.
- Carrillo, E., Matese, A., Rousseau, J., & Tisseyre, B. (2016). Use of multi-spectral airborne imagery to improve yield sampling in viticulture. *Precision Agriculture*, 17(1), 74-92.
- Oger, B., Roux, S., Le Moguédec, G., Tisseyre, B., (2020). Model and auxiliary data for an accurate estimate of the field mean. 13th European Conference on Precision Agriculture (ECPA 2021).
- Wasserman, L. (2004). All of statistics: A concise course in statistical inference. *New York: Springer*.