

# Revisiting Akaike's Final Prediction Error and the Generalized Cross Validation Criteria in Regression from the Same Perspective: From Least Squares to Ridge Regression and Smoothing Splines

Jean Raphael Ndzinga Mvondo\*, Eugène-Patrice Ndong Nguéma

Department of Mathematics, Ecole Polytechnique, Yaoundé I, Cameroon

Email: \*rndzinga@yahoo.fr

**How to cite this paper:** Ndzinga Mvondo, J.R. and Ndong Nguéma, E.-P. (2023) Revisiting Akaike's Final Prediction Error and the Generalized Cross Validation Criteria in Regression from the Same Perspective: From Least Squares to Ridge Regression and Smoothing Splines. *Open Journal of Statistics*, 13, 694-716.

<https://doi.org/10.4236/ojs.2023.135033>

**Received:** August 2, 2023

**Accepted:** September 25, 2023

**Published:** September 28, 2023

Copyright © 2023 by author(s) and Scientific Research Publishing Inc. This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

## Abstract

In regression, despite being both aimed at estimating the Mean Squared Prediction Error (MSPE), Akaike's Final Prediction Error (FPE) and the Generalized Cross Validation (GCV) selection criteria are usually derived from two quite different perspectives. Here, settling on the most commonly accepted definition of the MSPE as the expectation of the squared prediction error loss, we provide theoretical expressions for it, valid for any linear model (LM) fitter, be it under random or non random designs. Specializing these MSPE expressions for each of them, we are able to derive closed formulas of the MSPE for some of the most popular LM fitters: Ordinary Least Squares (OLS), with or without a full column rank design matrix; Ordinary and Generalized Ridge regression, the latter embedding smoothing splines fitting. For each of these LM fitters, we then deduce a computable estimate of the MSPE which turns out to coincide with Akaike's FPE. Using a slight variation, we similarly get a class of MSPE estimates coinciding with the classical GCV formula for those same LM fitters.

## Keywords

Linear Model, Mean Squared Prediction Error, Final Prediction Error, Generalized Cross Validation, Least Squares, Ridge Regression

## 1. Introduction

In many branches of activity, the data analyst is confronted with the need to model

a continuous numeric variable  $Y$  (the *response*) in terms of one or more explanatory other variables (called *covariates* or *predictors* or *regressors*)  $X_1, \dots, X_p$  in a population  $\Omega$  through a model

$$Y = \Phi(X) + \varepsilon, \tag{1}$$

where<sup>1</sup>:  $X = (X_1, \dots, X_p)^T$ ;  $\Phi(\cdot)$  is a function from  $\mathbb{R}^p \rightarrow \mathbb{R}$ , generally unknown, called the *regression function*;  $\varepsilon$  is an unobserved *error term*, also called *residual error* in the model (1).

In our developments in this paper, and contrary to popular tradition, we will not need that the variables  $X$  and  $\varepsilon$  necessarily be stochastically independent. However, the usual minimal assumption for (1) is that the variables  $X$  and  $\varepsilon$  satisfy:

$$(\mathcal{A}_0). \quad \mathbb{E}(\varepsilon | X) = 0.$$

Though sometimes far more debatable, we will also admit that the postulated regression model (1) satisfies the *homoscedasticity* assumption for the residual error variance:

$$(\mathcal{A}_1). \quad \text{Var}(\varepsilon | X) = \sigma^2 > 0.$$

Now, as is well known, Assumption  $(\mathcal{A}_0)$  implies that

$$\Phi(X) = \mathbb{E}(Y | X), \quad \text{i.e.} \quad \forall \mathbf{x} \in \mathbb{R}^p, \quad \Phi(\mathbf{x}) = \mathbb{E}(Y | X = \mathbf{x}). \tag{2}$$

However, that conditional expectation function can almost never be analytically computed in practical situations. The aim of the regression analysis of  $Y | X$  (*i.e.* “ $Y$  given  $X$ ”) is rather to estimate the unknown regression function  $\Phi(\cdot)$  in (1) based on some observed data

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^p \times \mathbb{R}, \quad \text{i.i.d.} \quad \overset{\mathcal{L}}{\sim} (X, Y), \tag{3}$$

collected on a sample of size  $n$  drawn from  $\Omega$ . If achieved, this will result in a computable function  $\hat{\Phi} : \mathbb{R}^p \rightarrow \mathbb{R}$  so that the final practical model used to express the response variable  $Y$  in terms of the vector of predictors  $X$  will be:

$$Y = \hat{\Phi}(X) + \tilde{\varepsilon}, \tag{4}$$

where  $\tilde{\varepsilon}$  is the residual error in the modeling. But once we get such a fit for the regression model (1), an obvious question then arises: how can one measure the accuracy of that computed fit (its so called *goodness-of-fit*)?

For some specific regression models (generally parametric ones and, most notably, the LM fitted by OLS with a full column rank design matrix), various measures of accuracy of their fit to given data, with closed form formulas, have been developed. But such specific and easily computable measures of accuracy are not universally applicable to all regression models. So they cannot be used to compare two arbitrarily fitted such models. Opposite to that, in the late 1960s, Akaike decisively introduced an approach (but in a much broader context in-

<sup>1</sup>By default, in this paper, vectors and random vectors are column vectors, unless specified otherwise; and  $A^T$  denotes the transpose of the matrix (or vector)  $A$ .

cluding time series modeling) which has that desirable universal feature [1], and is thus now generally recommended to compare different regression fits for the same data, albeit, sometimes, at some computational cost. It is based on estimating the *prediction error* of the fit. But various definitions and measures of that error have been used. By far, the most popular is the *Mean Squared Prediction Error* (MSPE), but there appears to be a myriad ways of defining it and/or theoretical estimates of it. A detailed lexicon on the matter is given in [2] and the references therein.

These various definitions and theoretical estimates of the MSPE are, undoubtedly, insightful in their respective motivations and aims, but they ultimately make the subject of prediction error assessment utterly complicated and confusing for most non expert users. This is compounded by the fact that many of these definitions and discussions of prediction error liberally mix theoretical estimates (*i.e.* finite averages over sample items) with intrinsic numeric characteristics of the whole population (expressed through *expectations* of some random variable). In that respect, while being both aimed at estimating the MSPE, the respective classical derivations of Akaike's Final Prediction Error (FPE) and Craven and Wahba's Generalized Cross Validation (GCV) stem from two quite different perspectives [1] [3]. This makes it not easy, for the non expert user, to grasp that these two selection criteria might even be related in the first place.

The first purpose of this paper is to settle on the definition of the MSPE most commonly known by users to assess the prediction power of any fitted model, be it for regression or else. Then, in that framework, we shall provide a conceptually simpler derivation of the FPE as an estimate of that MSPE in a LM fitted by OLS when the design matrix is full column rank. Secondly, we build on that to derive generalizations of the FPE for the LM fitted by other well known methods under various scenarios, generalizations seldom accessible from the traditional derivation of the FPE. Finally, we show that, in that same unified framework, a minor variation in the derivation of the MSPE estimates yield the well known formula of the GCV for all these various LM fitters. For the latter selection criterion, previous attempts [4] [5] have been made to provide a derivation of it different from the classical one as an approximation of the leave-one-out cross validation (LOO-CV) estimate of the MSPE. We view our approach as a more straightforward and inclusive derivation of the GCV score.

To achieve that, we start, in Section 2, by reviewing the prediction error viewpoint in assessing how a fitted regression model performs and to settle on the definition of the MSPE most commonly known by users to measure that performance (while briefly recalling the alternative one most given in regression textbooks). Then, in Section 3, focusing specifically on the LM, we provide theoretical expressions of that MSPE measure valid for any arbitrary LM fitting method, be it under random or non random designs. In the next sections, these expressions are successively specialized to some of the best known LM fitters to deduce, for each, computable estimates of the MSPE, a class of which yielding

the FPE and a slight variation the GCV. Those LM fitters include: OLS, with or without a full column rank design matrix (Section 4), Ridge Regression, both Ordinary and Generalized (Section 5), the latter embedding smoothing splines fitting. Finally, Section 6 draws a conclusion and suggests some needed future work.

As customary, we summarize the data (3) through the design matrix and the response vector:

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathcal{M}_{n,p}(\mathbb{R}) \quad \text{and} \quad \mathbf{y} = (y_1, \dots, y_n)^T \in \mathbb{R}^n. \quad (5)$$

It is important to emphasize that technically, each observed response  $y_i$  is a realization of a random variable. Consequently, the same will be true of the vector  $\mathbf{y}$ . On the other hand, the matrix  $\mathbf{X}$  is considered as fixed or random, depending on whether the covariates  $X_1, \dots, X_p$  are viewed as fixed or random variables. According to which one holds, one talks of *fixed design* or *random design*. Our presentation will encompass both scenarios. For convenience, we will use the same symbol to denote a random variable and its realization. Nonetheless, the distinction between the latter two will be apparent from context.

## 2. Prediction Error of a Fitted Regression Model

### 2.1. The Prediction Error Viewpoint for Assessing a Fitted Regression Model

Using the data (3), assume we got an estimate  $\hat{\Phi}(\cdot)$  of the unknown regression function  $\Phi(\cdot)$  in the regression model (1). Thus, we fitted the model (1) to express the response  $Y$  in terms of the vector of covariates  $X = (X_1, \dots, X_p)$  in the population under study  $\Omega$ . From a predictive perspective, assessing the goodness-of-fit of that fitted model amounts to answering the question: *How well the fitted model is likely to predict the response variable  $Y$  on a future individual based on its covariates values  $X = (X_1, \dots, X_p)$ ?*

To try to formalize an answer, consider a new member of  $\Omega$ , drawn independently from the sample of  $\Omega$  which produced the data (3), and assume we have observed its covariates vector  $X = X_0 \in \mathbb{R}^p$ , but not its response value  $Y = Y_0 \in \mathbb{R}$ . Nonetheless, served value of the latter. This is the *prediction problem of  $Y|X$*  (“ $Y$  given  $X$ ”) on that individual. With model (4) fitted to the response variable  $Y$ , it seems natural to predict the unknown value of  $Y = Y_0$  on that individual, given from the former, we want to get an idea of the unob  $X = X_0$ , by:

$$\hat{Y}_0 = \hat{\Phi}(X_0). \quad (6)$$

But there is, obviously, an error attached to such a prediction of  $Y = Y_0$  value by  $\hat{Y}_0$ , called the *prediction error* or *predictive risk* of  $\hat{Y}_0$  w.r.t.  $Y_0$ .

### 2.2. The Mean Squared Prediction Error

The prediction error of  $\hat{Y}_0$  w.r.t.  $Y_0$  needs to be assessed beforehand to get an idea of the quality of the fit provided by the estimated regression model (4) for

$Y|X$  in the population under study. To that aim, the global measure mostly used is the Mean Squared Prediction Error (MSPE), but one meets a bewildering number of ways of defining it in the literature [2], creating a bit of confusion in the mindset of the daily user of Statistics. Yet, most users would accept as definition of the MSPE:

$$\text{MSPE}(\hat{Y}_0, Y_0) = \mathbb{E}(\hat{Y}_0 - Y_0)^2, \tag{7}$$

an intrinsic characteristic of the population which we need to estimate from the available data (3).

In this paper, an unconditional expectation like the *r.h.s.* of (7) is meant to integrate over all the random variables involved in the integrand. A conditional expectation, on the other hand, acts the same, except for the conditioning random variable(s). As such, *en route* to getting  $\text{MSPE}(\hat{Y}_0, Y_0)$ , we will pass successively through:

$$\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}, \mathbf{y}, X_0) = \mathbb{E} \left[ (\hat{Y}_0 - Y_0)^2 | \mathbf{X}, \mathbf{y}, X_0 \right], \tag{8a}$$

$$\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}, \mathbf{y}) = \mathbb{E} \left[ (\hat{Y}_0 - Y_0)^2 | \mathbf{X}, \mathbf{y} \right], \tag{8b}$$

$$\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}) = \mathbb{E} \left[ (\hat{Y}_0 - Y_0)^2 | \mathbf{X} \right], \tag{8c}$$

representing, each, the MSPE conditioned on some information, which might be relevant in its own right. In particular,  $\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X})$  is the relevant definition of the MSPE in the case of a fixed design. But, as a consequence of moving successively through (8a)-(8c) to get (7), handling the fixed design case will not need a special treatment because computing (8c) will be a necessary intermediate step.

Obviously, for most fitted regression models, trying to compute  $\text{MSPE}(\hat{Y}_0, Y_0)$  or  $\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X})$  is a hopeless task. However, it is known that  $K$ -fold cross validation (CV) can be used to estimate these quantities in a nearly universal manner, imposing no distributional assumption and using a quasi automatic algorithm [6]. Nonetheless, cross validation has its own defects such as a high extra computational cost<sup>2</sup>, the impact of the choice of  $K$  and, generally, upwards bias. As for the latter defect, Borra and Di Ciaccio [2] showed in an extensive simulation study that Repeated Corrected CV, a little popularized correction to  $K$ -fold cross validation developed by Burman [7], performed pretty well and outperformed, on some well known nonlinear regression fitters and under a variety of scenarios, all the other theoretically more technically involved suggested measures of the MSPE alluded to above.

Fortunately, after fitting a Linear Model to given data by a chosen method,

<sup>2</sup>However, that defect is more and more mitigated these days, thanks to the availability of increasingly user-friendly parallel programming environments in popular statistical software systems, provided one has a computer allowing to exploit such possibilities, *e.g.* a laptop with several *core processors*.

there is, at least for the most common methods, no need to shoulder the computational cost attached to CV to estimate the MSPE of the fitted model. It is the purpose of this article to show that for the most well known LM fitters, one can transform (7)-(8c) in a form allowing to deduce, in a straightforward manner, two estimates of the MSPE computable from the data, which turn out to coincide, respectively, with the FPE and the GCV. This, therefore, will yield a new way to get these two selection criteria in linear regression completely different from how they are usually respectively motivated and derived.

### 2.3. Prediction Error and Sources of Randomness in the Regression Process

The rationale behind settling on (7) or (8c) as definition of the MSPE lies in the fact that, from our standpoint, any global measure of the prediction error in the population computed from a sample collected in it must account both for the randomness in drawing that sample and in that of drawing a future individual. Consequently, each of the expectations in (7)-(8c) should be understood as integrating over all the possible sources of randomness entering in the process of computing the prediction  $\hat{Y}_0$  of  $Y_0$ , save the conditioning variables, if any:

1) the realized, but unobserved random  $n$ -vector  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top$  of sample residual errors in the model (1) for the data (3). This vector results from the fact that for the observed sample of  $n$  individuals, the model (1) implies that

$$y_i = \Phi(\mathbf{x}_i) + \varepsilon_i, \quad i = 1, \dots, n, \quad \text{with } \varepsilon_1, \dots, \varepsilon_n \text{ i.i.d. } \stackrel{\mathcal{L}}{\sim} \varepsilon; \quad (9)$$

2)  $X_0 \in \mathbb{R}^p$ , the vector of covariates for the potential newly independently drawn individual for which the unknown value  $Y_0 \in \mathbb{R}$  of the response  $Y$  is to be predicted;

3)  $\varepsilon_0$ , the error of the model (1) for that new individual,  $Y_0 = \Phi(X_0) + \varepsilon_0$ ,  $\varepsilon_0 \stackrel{\mathcal{L}}{\sim} \varepsilon$ .

4) and, in case of a random design, the entries in the design matrix  $X$ .

The key assumption to assess the prediction error of a regression model is then:

( $\mathcal{A}_2$ ). The random couple  $(\varepsilon_0, X_0) \stackrel{\mathcal{L}}{\sim} (\varepsilon, X)$  and is independent from  $(\varepsilon, X)$ .

### 2.4. Measuring the Prediction Error in Regression: Sample Based Definition

While the GCV score is classically derived, indeed, as an estimate of the MSPE as given by (7), through a two-stage process where the intermediate step is the well known LOO-CV estimate of that MSPE, the traditional derivation of the FPE criterion stems from a completely different angle. Actually, the latter angle is the one presented in most textbooks on regression [8] [9]. In it, with the regression function  $\Phi$  in (1) estimated through  $\hat{\Phi}$ , a function computed from the data (3), the prediction error of that fit is rather measured by how well the vector of predicted responses on sample items,  $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)^\top$ , with  $\hat{y}_i = \hat{\Phi}(\mathbf{x}_i)$ ,  $i = 1, \dots, n$ ,

estimates the vector of exact responses on those items,  $\mathbf{\Phi} = (\Phi_1, \dots, \Phi_n)^T$ , with  $\Phi_i = \Phi(\mathbf{x}_i)$ ,  $i = 1, \dots, n$ . In that respect, one considers the *mean average squared error*, also called *risk*,

$$\text{MASE}(\hat{\mathbf{y}} | \mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(\hat{y}_i - \Phi_i)^2 | \mathbf{X}], \tag{10}$$

with bias-variance decomposition

$$\text{MASE}(\hat{\mathbf{y}} | \mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \text{Var}(\hat{y}_i | \mathbf{X}) + \frac{1}{n} \sum_{i=1}^n b_i^2. \tag{11}$$

where  $b_i = \mathbb{E}(\hat{y}_i | \mathbf{X}) - \Phi_i$  is the conditional bias, given  $\mathbf{X}$ , in estimating  $\Phi_i$  by  $\hat{y}_i$ .

But the relation to the more obvious definition (8c) of the conditional MSPE is better seen when one considers, instead, the *average predicted squared error* ([8] Chapter 3, page 42) or *prediction risk* ([9] Chapter 2, page 29):

$$\text{PSE}(\hat{\mathbf{y}} | \mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y_i^* - \hat{y}_i)^2 | \mathbf{X}], \tag{12}$$

where  $y_1^*, \dots, y_n^*$  are putative responses assumed generated at the respective predictor values  $\mathbf{x}_1, \dots, \mathbf{x}_n$  through model (1), but with respective errors  $\varepsilon_1^*, \dots, \varepsilon_n^*$  independent from the initial ones  $\varepsilon_1, \dots, \varepsilon_n$ . Nonetheless, there is a simple relation between (10) and (12):

$$\text{PSE}(\hat{\mathbf{y}} | \mathbf{X}) = \sigma^2 + \text{MASE}(\hat{\mathbf{y}} | \mathbf{X}). \tag{13}$$

hence minimizing  $\text{PSE}(\hat{\mathbf{y}} | \mathbf{X})$  w.r.t.  $\hat{\mathbf{\Phi}}$  is the same as doing so for  $\text{MASE}(\hat{\mathbf{y}} | \mathbf{X})$ .

In its classical derivation for linear regression (see, e.g., [10], pages 19-20), the FPE selection criterion is an estimate of  $\text{PSE}(\hat{\mathbf{y}} | \mathbf{X})$ . Now, with the terminology elaborated in [5], measuring the prediction error by the latter amounts to using a *Fixed-X* viewpoint as opposed to the *Random-X* one when measuring it instead through  $\text{MSPE}(\hat{Y}_0, Y_0)$ . But an even more common terminology to distinguish these two approaches to estimating the predictive performance of a regression method qualifies the first as *in-sample* prediction and the second one as *out-of-sample* prediction. It should be said that while in the past, the prediction error was mostly evaluated using the in-sample paradigm, facing the complexity of data met in modern statistics, noticeably high dimensional data, many researchers in regression have advocated or used the out-of-sample viewpoint, though this might be through either (7), (8b), or (8c), depending on the author(s). In that respect, in addition to the aforementioned paper of Trosset and Tibshirani, we may cite, e.g., Breiman and Spector [11], Leeb [12], Dicker [13], Dobriban and Wager [14].

Note, however, that the prediction error viewpoint in assessing the quality of a regression fit is not without its own demerits. Indeed, in [15] and [16], it is highlighted that in the specific case of smoothing splines regression, one can find a fit which is optimal from the prediction error viewpoint, but which clearly un-

dersmooths the data, resulting in a wiggly curve. But the argument appears to be more a matter of visual esthetics because the analysis in those papers targets the regression function  $\Phi$ , which is, indeed, probably the main objective of many users of univariate nonparametric regression. Nonetheless, when measuring the prediction error through  $\text{MSPE}(\hat{Y}_0, Y_0)$ , the target is rather the response  $Y$ , *i.e.*  $\Phi + \text{error}$ , in which the wiggleness is inherently embedded. Do not forget that when formulating his Final Prediction Error, Akaike was working on real world engineering problems [17]. Hence his interest in targeting the actual output in his predictions.

### 3. MSPE of a Linear Model Fitter

From now on, we focus attention on the generic LM, with  $\beta = (\beta_1, \dots, \beta_p)^\top$  :

$$Y = \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon = X^\top \beta + \varepsilon, \tag{14}$$

to be fitted to the data (3). Because of its ease of tractability and manipulation, the LM is the most popular approach to estimating the regression function  $\Phi(\cdot)$  in (1). It is mostly implemented by estimating  $\beta$  through the Ordinary Least Squares criterion. However, several other approaches have been designed to estimate  $\beta$  based on various grounds, such as Weighted Least Squares, Total Least Squares, Least Absolute Deviations, LASSO [18] and Ridge Regression [19]. Furthermore, some generally more adaptive nonparametric regression methods proceed by first nonlinearly transforming the data to a scale where they can be fitted by a LM. Due to its more than ever quite central role in statistical modeling and practice, several books have been and continue to be fully devoted to the presentation of the LM and its many facets such as: [20] [21] [22] and [23].

#### 3.1. Some Preliminaries

For the sample of  $n$  individuals with recorded data (3), the general regression Equation (9) becomes, in the case of the LM (14):

$$y_i = \mathbf{x}_i^\top \beta + \varepsilon_i, \quad i = 1, \dots, n, \tag{15}$$

or, better globally summarized in matrix form,

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon, \quad \text{with } \varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top. \tag{16}$$

Then Assumptions  $(\mathcal{A}_0)$  and  $(\mathcal{A}_1)$  respectively imply here:

$$\mathbb{E}(\varepsilon | \mathbf{X}) = \mathbf{0} \quad \text{and} \quad \mathbb{M}\text{cov}(\varepsilon | \mathbf{X}) = \sigma^2 \mathbf{I}_n, \tag{17}$$

with  $\mathbf{I}_n$  the  $n$ -by- $n$  identity matrix. Consequently,

$$\mathbb{E}(\mathbf{y} | \mathbf{X}) = \mathbf{X}\beta \quad \text{and} \quad \mathbb{M}\text{cov}(\mathbf{y} | \mathbf{X}) = \sigma^2 \mathbf{I}_n. \tag{18}$$

Fitting the LM (14) boils down to estimating the  $p$ -vector of parameters  $\beta$ . We call *Linear Model fitter*, or *LM fitter*, any method allowing to achieve that. First, we consider an arbitrarily chosen such method. It uses the data (3) to



compute  $\hat{\beta} \in \mathbb{R}^p$ , an estimate of  $\beta$ . The precision of that estimate<sup>3</sup> can be assessed through its *Mean Squared Error matrix*:

$$\text{Mse}(\hat{\beta}) = \mathbb{E}\left[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^\top\right]. \tag{19a}$$

But to reach  $\text{Mse}(\hat{\beta})$  generally requires passing through the conditional Mean Squared Error matrix of  $\hat{\beta}$  given the design matrix  $\mathbf{X}$ :

$$\text{Mse}(\hat{\beta} | \mathbf{X}) = \mathbb{E}\left[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^\top | \mathbf{X}\right]. \tag{19b}$$

The relationship between the two is:

$$\text{Mse}(\hat{\beta}) = \mathbb{E}_{\mathbf{X}}\left[\text{Mse}(\hat{\beta} | \mathbf{X})\right]. \tag{19c}$$

Those two matrices will play a key role in our MSPE derivations to come.

The precision of an estimate of a vector parameter like  $\beta$  is easier to assess when its Mean Squared Error matrix coincides with its covariance matrix. Hence, our interest in:

**Definition 1** In the LM (14),  $\hat{\beta}$  is an unbiased estimate of  $\beta$ , conditionally on  $\mathbf{X}$ , if  $\mathbb{E}(\hat{\beta} | \mathbf{X}) = \beta$ .

Then one has:  $\text{Mse}(\hat{\beta} | \mathbf{X}) = \text{Mcov}(\hat{\beta} | \mathbf{X})$ , the conditional covariance matrix of  $\hat{\beta}$  given  $\mathbf{X}$ .

Since  $\mathbb{E}(\hat{\beta}) = \mathbb{E}_{\mathbf{X}}\left[\mathbb{E}(\hat{\beta} | \mathbf{X})\right]$ , it is immediate that if  $\hat{\beta}$  is an unbiased estimate of  $\beta$ , conditionally on  $\mathbf{X}$ , then  $\mathbb{E}(\hat{\beta}) = \beta$ , i.e.  $\hat{\beta}$  is an unbiased estimate of  $\beta$  (unconditionally). So the former property is stronger than the latter, but is more useful in this context. Note also that it implies:  $\text{Mse}(\hat{\beta}) = \text{Mcov}(\hat{\beta})$ , the covariance matrix of  $\hat{\beta}$ . On the other hand, when  $\hat{\beta}$  is a biased estimate of  $\beta$ , the bias-variance decomposition of  $\text{Mse}(\hat{\beta})$  might be of interest:

$$\text{Mse}(\hat{\beta}) = \text{Bias}(\hat{\beta})\text{Bias}(\hat{\beta})^\top + \text{Mcov}(\hat{\beta}), \tag{20}$$

where  $\text{Bias}(\hat{\beta}) = \mathbb{E}(\hat{\beta}) - \beta$ .

### 3.2. MSPE of a Linear Model Fitter: Base Results

In the prediction setting of Section 2.1 applied to the LM (14), one has, for the new individual:

$$Y_0 = X_0^\top \beta + \varepsilon_0, \tag{21}$$

with  $\varepsilon_0 \stackrel{\mathcal{L}}{\sim} \varepsilon$ , but unknown. It would then be natural to predict the response value  $Y_0$  by

$$\tilde{Y}_0 = X_0^\top \beta = \Phi(X_0), \tag{22}$$

were the exact value of the parameter vector  $\beta$  available. Since that is not typ-

<sup>3</sup>Keeping in line with our stated convention of denoting a random variable and its realization by the same symbol, whereas, technically, an *estimate* of a parameter is a realization of a random variable called *estimator* of that parameter, we use the term *estimate* here for both. So when an estimate appears inside an expectation or a covariance notation, it is definitely the estimator which is meant.

ically the case, one rather predicts  $Y_0$  by the computable quantity

$$\hat{Y}_0 = X_0^T \hat{\beta} = \hat{\Phi}(X_0). \tag{23}$$

The goal here is to find expressions, in this context, for the MSPE of  $\hat{Y}_0$  w.r.t.  $Y_0$  as given by (7)-(8c), manageable enough to allow the derivation of computable estimates of that MSPE for the most common LM fitters. The starting point to get such MSPE expressions is the base result:

**Theorem 1** For any  $\hat{\beta}$  estimating  $\beta$  in the LM (14), one has, under  $(\mathcal{A}_0)$ ,  $(\mathcal{A}_1)$ ,  $(\mathcal{A}_2)$ :

$$\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}, \mathbf{y}, X_0) = \sigma^2 + \text{tr} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)^T \cdot X_0 X_0^T \right], \tag{24a}$$

$$\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}, \mathbf{y}) = \sigma^2 + \text{tr} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)^T \cdot \mathbb{E}(X_0 X_0^T) \right], \tag{24b}$$

$$\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}) = \sigma^2 + \text{tr} \left[ \text{Mlse}(\hat{\beta} | \mathbf{X}) \cdot \mathbb{E}(X_0 X_0^T) \right], \tag{24c}$$

$$\text{MSPE}(\hat{Y}_0, Y_0) = \sigma^2 + \text{tr} \left[ \text{Mlse}(\hat{\beta}) \cdot \mathbb{E}(X_0 X_0^T) \right]. \tag{24d}$$

*Proof.* From  $\hat{Y}_0 - Y_0 = X_0^T (\hat{\beta} - \beta) - \varepsilon_0$ , we first get:

$$(\hat{Y}_0 - Y_0)^2 = \varepsilon_0^2 - 2\varepsilon_0 X_0^T (\hat{\beta} - \beta) + \left[ X_0^T (\hat{\beta} - \beta) \right]^2. \tag{25}$$

Now, from (8a) and using Assumption  $(\mathcal{A}_2)$ ,

$$\begin{aligned} \text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}, \mathbf{y}, X_0) &= \mathbb{E}(\varepsilon_0^2 | X_0) - 2\mathbb{E}(\varepsilon_0 | X_0) X_0^T (\hat{\beta} - \beta) + \left[ X_0^T (\hat{\beta} - \beta) \right]^2 \\ &= \text{Var}(\varepsilon_0 | X_0) + \left[ X_0^T (\hat{\beta} - \beta) \right]^2 = \sigma^2 + \left[ X_0^T (\hat{\beta} - \beta) \right]^2, \end{aligned} \tag{26}$$

the latter because  $\mathbb{E}(\varepsilon_0 | X_0) = \mathbb{E}(\varepsilon | X) = 0$ , and so

$\mathbb{E}(\varepsilon_0^2 | X_0) = \text{Var}(\varepsilon_0 | X_0) = \text{Var}(\varepsilon | X) = \sigma^2$ . On the other hand,

$\left[ X_0^T (\hat{\beta} - \beta) \right]^2 = \text{tr} \left[ X_0^T (\hat{\beta} - \beta) \right]^2$  because  $\left[ X_0^T (\hat{\beta} - \beta) \right]^2$  is a scalar. Thus

$$\left[ X_0^T (\hat{\beta} - \beta) \right]^2 = \text{tr} \left[ X_0^T (\hat{\beta} - \beta) (\hat{\beta} - \beta)^T X_0 \right] = \text{tr} \left[ (\hat{\beta} - \beta) (\hat{\beta} - \beta)^T X_0 X_0^T \right],$$

which, inserted in (26), yields (24a).

Thanks to (24a) and identity (8b), we have

$$\begin{aligned} \text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}, \mathbf{y}) &= \mathbb{E} \left[ (\hat{Y}_0 - Y_0)^2 | \mathbf{X}, \mathbf{y} \right] \\ &= \mathbb{E}_{X_0} \left\{ \mathbb{E} \left[ (\hat{Y}_0 - Y_0)^2 | \mathbf{X}, \mathbf{y}, X_0 \right] \right\}, \text{ since } X_0 \perp\!\!\!\perp (\mathbf{X}, \mathbf{y}) \\ &= \mathbb{E}_{X_0} \left\{ \sigma^2 + \text{tr} \left[ (\hat{\beta} - \beta) (\hat{\beta} - \beta)^T \cdot X_0 X_0^T \right] \right\} \\ &= \mathbb{E}_{X_0} (\sigma^2) + \mathbb{E}_{X_0} \left\{ \text{tr} \left[ (\hat{\beta} - \beta) (\hat{\beta} - \beta)^T \cdot X_0 X_0^T \right] \right\} \end{aligned}$$

$$\begin{aligned}
 &= \sigma^2 + \text{tr} \left\{ \mathbb{E}_{X_0} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)^T \cdot X_0 X_0^T \right] \right\} \\
 &= \sigma^2 + \text{tr} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)^T \cdot \mathbb{E}_{X_0} (X_0 X_0^T) \right] \\
 &= \sigma^2 + \text{tr} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)^T \cdot \mathbb{E} (X_0 X_0^T) \right].
 \end{aligned}$$

Likewise, (24b) and (8c) give:

$$\begin{aligned}
 \text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}) &= \mathbb{E} \left[ (\hat{Y}_0 - Y_0)^2 | \mathbf{X} \right] = \mathbb{E}_{y|\mathbf{X}} \mathbb{E} \left[ (\hat{Y}_0 - Y_0)^2 | \mathbf{X}, \mathbf{y} \right] \\
 &= \mathbb{E}_{y|\mathbf{X}} \left\{ \sigma^2 + \text{tr} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)^T \cdot \mathbb{E} (X_0 X_0^T) \right] \right\} \\
 &= \mathbb{E}_{y|\mathbf{X}} (\sigma^2) + \mathbb{E}_{y|\mathbf{X}} \left\{ \text{tr} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)^T \cdot \mathbb{E} (X_0 X_0^T) \right] \right\} \\
 &= \sigma^2 + \mathbb{E}_{y|\mathbf{X}} \left\{ \text{tr} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)^T \cdot \mathbb{E} (X_0 X_0^T) \right] \right\} \\
 &= \sigma^2 + \text{tr} \left\{ \mathbb{E}_{y|\mathbf{X}} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)^T \cdot \mathbb{E} (X_0 X_0^T) \right] \right\} \\
 &= \sigma^2 + \text{tr} \left\{ \mathbb{E}_{y|\mathbf{X}} \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)^T \right] \cdot \mathbb{E} (X_0 X_0^T) \right\} \\
 &= \sigma^2 + \text{tr} \left[ \text{Mise}(\hat{\beta} | \mathbf{X}) \cdot \mathbb{E} (X_0 X_0^T) \right].
 \end{aligned}$$

From relations (24c) and (7), one gets:

$$\begin{aligned}
 \text{MSPE}(\hat{Y}_0, Y_0) &= \mathbb{E} (\hat{Y}_0 - Y_0)^2 = \mathbb{E}_X \mathbb{E} \left[ (\hat{Y}_0 - Y_0)^2 | \mathbf{X} \right] \\
 &= \mathbb{E}_X \left\{ \sigma^2 + \text{tr} \left[ \text{Mise}(\hat{\beta} | \mathbf{X}) \cdot \mathbb{E} (X_0 X_0^T) \right] \right\} \\
 &= \mathbb{E}_X (\sigma^2) + \mathbb{E}_X \left\{ \text{tr} \left[ \text{Mise}(\hat{\beta} | \mathbf{X}) \cdot \mathbb{E} (X_0 X_0^T) \right] \right\} \\
 &= \sigma^2 + \text{tr} \left\{ \mathbb{E}_X \left[ \text{Mise}(\hat{\beta} | \mathbf{X}) \cdot \mathbb{E} (X_0 X_0^T) \right] \right\} \\
 &= \sigma^2 + \text{tr} \left\{ \mathbb{E}_X \left[ \text{Mise}(\hat{\beta} | \mathbf{X}) \right] \cdot \mathbb{E} (X_0 X_0^T) \right\} \\
 &= \sigma^2 + \text{tr} \left[ \text{Mise}(\hat{\beta}) \cdot \mathbb{E} (X_0 X_0^T) \right].
 \end{aligned}$$

□

The above result is interesting in that it imposes no assumption on  $\hat{\beta}$ , hence it is valid for any LM fitter. But an immediate important subcase is provided in:

**Corollary 2** *If, conditional on  $\mathbf{X}$ ,  $\hat{\beta}$  estimates  $\beta$  unbiasedly in the LM (14), then, under Assumptions  $(\mathcal{A}_0)$ ,  $(\mathcal{A}_1)$  and  $(\mathcal{A}_2)$ :*

$$\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}) = \sigma^2 + \text{tr} \left[ \text{Mlcov}(\hat{\beta} | \mathbf{X}) \cdot \mathbb{E} (X_0 X_0^T) \right], \tag{27a}$$

$$\text{MSPE}(\hat{Y}_0, Y_0) = \sigma^2 + \text{tr} \left[ \text{Mlcov}(\hat{\beta}) \cdot \mathbb{E} (X_0 X_0^T) \right]. \tag{27b}$$

#### 4. MSPE When Fitting the LM by Ordinary Least Squares

By far, the most popular approach to estimating the parameter vector  $\beta$  in the LM (14) is through minimizing the Ordinary Least Squares (OLS) criterion us-

ing the observed data (3):

$$\hat{\beta} = \hat{\beta}_{OLS} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2. \tag{28}$$

The properties of  $\hat{\beta}_{OLS}$  as an estimate of  $\beta$  depend on whether the design matrix  $\mathbf{X}$  has full column rank or not. This remains true when studying the corresponding MSPE as well.

### 4.1. MSPE in the LM Fitted by OLS with $\mathbf{X}$ of Full Column Rank

#### 4.1.1. LM Fitted by OLS with $\mathbf{X}$ Full Column Rank

Here, we consider the LM (14) fitted under:

( $\mathcal{A}_3$ ).  $\operatorname{rank}(\mathbf{X}) = p$ , i.e.  $\mathbf{X}$  is a full column rank matrix.

That assumption is known to be equivalent to saying that the square matrix  $\mathbf{X}^T \mathbf{X}$  is nonsingular, thus implying that

$$\hat{\beta}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}. \tag{29}$$

Furthermore, given Assumptions ( $\mathcal{A}_0$ )-( $\mathcal{A}_1$ ), (18) holds, so

$$\mathbb{E}(\hat{\beta}_{OLS} | \mathbf{X}) = \beta \quad \text{and} \quad \mathbb{M}\operatorname{cov}(\hat{\beta}_{OLS} | \mathbf{X}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}. \tag{30}$$

We also recall that under these assumptions, with  $\hat{\mathbf{e}} = \mathbf{y} - \mathbf{X} \hat{\beta}_{OLS}$  the residual response vector and  $\|\cdot\|$  the Euclidean norm, a computable unbiased estimate of the residual variance  $\sigma^2$  is:

$$\hat{\sigma}^2 = \frac{\operatorname{SSR}}{n-p}, \quad \text{with} \quad \operatorname{SSR} = \|\hat{\mathbf{e}}\|^2 = \sum_{i=1}^n (y_i - \mathbf{x}_i^T \hat{\beta}_{OLS})^2, \tag{31}$$

the latter being the *sum of squared residuals* in the OLS fit to the data.

Now, from the first identity in (30), we deduce that when  $\mathbf{X}$  is full column rank,  $\hat{\beta}_{OLS}$  is an unbiased estimate of  $\beta$ , conditionally on  $\mathbf{X}$ . Hence, combining the second identity in (30) with Corollary 2 yields:

**Theorem 3** *In the LM (14) fitted by OLS with Assumptions ( $\mathcal{A}_0$ ), ( $\mathcal{A}_1$ ), ( $\mathcal{A}_2$ ) and ( $\mathcal{A}_3$ ),*

$$\operatorname{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}) = \sigma^2 \left( 1 + \operatorname{tr} \left[ (\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbb{E}(X_0 X_0^T) \right] \right), \tag{32a}$$

$$\operatorname{MSPE}(\hat{Y}_0, Y_0) = \sigma^2 \left( 1 + \operatorname{tr} \left\{ \mathbb{E} \left[ (\mathbf{X}^T \mathbf{X})^{-1} \right] \cdot \mathbb{E}(X_0 X_0^T) \right\} \right). \tag{32b}$$

*Proof.* From (24c),

$$\begin{aligned} \operatorname{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}) &= \sigma^2 + \operatorname{tr} \left[ \mathbb{M}\operatorname{se}(\hat{\beta}_{OLS} | \mathbf{X}) \cdot \mathbb{E}(X_0 X_0^T) \right] \\ &= \sigma^2 + \operatorname{tr} \left[ \mathbb{M}\operatorname{cov}(\hat{\beta}_{OLS} | \mathbf{X}) \cdot \mathbb{E}(X_0 X_0^T) \right] \\ &= \sigma^2 + \operatorname{tr} \left[ \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbb{E}(X_0 X_0^T) \right] \\ &= \sigma^2 + \sigma^2 \operatorname{tr} \left[ (\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbb{E}(X_0 X_0^T) \right] \\ &= \sigma^2 \left( 1 + \operatorname{tr} \left[ (\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbb{E}(X_0 X_0^T) \right] \right). \end{aligned}$$

Now, using (32a), one gets:

$$\begin{aligned}
 \text{MSPE}(\hat{Y}_0, Y_0) &= \mathbb{E}_X \mathbb{E} \left[ (\hat{Y}_0 - Y_0)^2 \mid \mathbf{X} \right] \\
 &= \mathbb{E}_X \left\{ \sigma^2 \left( 1 + \text{tr} \left[ (\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbb{E}(X_0 X_0^T) \right] \right) \right\} \\
 &= \sigma^2 \left\{ \mathbb{E}_X \left( 1 + \text{tr} \left[ (\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbb{E}(X_0 X_0^T) \right] \right) \right\} \\
 &= \sigma^2 \left\{ 1 + \mathbb{E}_X \left( \text{tr} \left[ (\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbb{E}(X_0 X_0^T) \right] \right) \right\} \\
 &= \sigma^2 \left\{ 1 + \text{tr} \left( \mathbb{E}_X \left[ (\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbb{E}(X_0 X_0^T) \right] \right) \right\} \\
 &= \sigma^2 \left( 1 + \text{tr} \left\{ \mathbb{E} \left[ (\mathbf{X}^T \mathbf{X})^{-1} \right] \cdot \mathbb{E}(X_0 X_0^T) \right\} \right),
 \end{aligned}$$

from which the result is got.  $\square$

#### 4.1.2. The FPE and the GCV in the LM Fitted by OLS with X of Full Column Rank

From (32b) in Theorem 3, we deduce a closed form computable estimate of  $\text{MSPE}(\hat{Y}_0, Y_0)$ , using data (3), by estimating, respectively:

- the residual variance  $\sigma^2$  by  $\hat{\sigma}^2$  given by (31);
- the  $p \times p$  expectation matrix  $\mathbb{E} \left[ (\mathbf{X}^T \mathbf{X})^{-1} \right]$  by the observed  $(\mathbf{X}^T \mathbf{X})^{-1}$ ;
- the  $p \times p$  expectation matrix  $\mathbb{E}(X_0 X_0^T)$  by (given that  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are *i.i.d.*  $\stackrel{\mathcal{L}}{\sim} X_0$ ):

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \frac{1}{n} \mathbf{X}^T \mathbf{X}. \tag{33}$$

Therefore, one estimates  $\text{MSPE}(\hat{Y}_0, Y_0)$  by

$$\begin{aligned}
 \widehat{\text{MSPE}}(\hat{Y}_0, Y_0) &= \hat{\sigma}^2 \left\{ 1 + \frac{1}{n} \text{tr} \left[ (\mathbf{X}^T \mathbf{X})^{-1} \cdot (\mathbf{X}^T \mathbf{X}) \right] \right\} = \hat{\sigma}^2 \left[ 1 + \frac{1}{n} \text{tr}(\mathbf{I}_p) \right] \\
 &= \hat{\sigma}^2 \left( 1 + \frac{p}{n} \right) = \frac{n+p}{n-p} \cdot \frac{\text{SSR}}{n} = \hat{S}^2 \cdot \frac{n+p}{n-p} = \text{FPE},
 \end{aligned} \tag{34}$$

with

$$\hat{S}^2 = \text{SSR}/n = \|\hat{\boldsymbol{\epsilon}}\|^2/n = \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_{\text{OLS}}\|^2/n, \tag{35}$$

the usual Maximum Likelihood Estimator of the residual variance  $\sigma^2$  in the LM (14) when the residual error  $\boldsymbol{\epsilon}$  is assumed to follow a  $\mathcal{N}(0, \sigma^2)$  Gaussian distribution.

We see that the final estimate  $\widehat{\text{MSPE}}(\hat{Y}_0, Y_0)$  obtained for  $\text{MSPE}(\hat{Y}_0, Y_0)$  coincides with Akaike’s Final Prediction Error (FPE) goodness-of-fit criterion for the LM (14) fitted by OLS [1]. The main difference between the derivation above and the traditional one is that the latter uses the sample viewpoint of the prediction error reviewed in Section 2.4. The latter excludes the possibility that covariates values on a future individual might be completely unrelated to the observed  $\mathbf{x}_i$ ’s in the sample (3). In particular, it does not account for any po-

tential random origin for the design matrix  $\mathbf{X}$ , a situation often encountered in certain areas of application of the LM such as in econometrics.

On the other hand, estimating, instead,  $\mathbb{E}(X_0 X_0^T)$  by  $\mathbf{X}^T \mathbf{X} / (n - p)$  yields as estimate of  $\text{MSPE}(\hat{Y}_0, Y_0)$ :

$$\hat{\sigma}^2 \left( 1 + \frac{p}{n - p} \right) = \frac{n \cdot \text{SSR}}{(n - p)^2} = \frac{\hat{S}^2}{(1 - p/n)^2} = \text{GCV}, \tag{36}$$

the traditional GCV estimate of  $\text{MSPE}(\hat{Y}_0, Y_0)$  in the LM fitted by OLS with  $\mathbf{X}$  full column rank.

**Remark 1** *Note that the very way the two estimates (34) and (36) of  $\text{MSPE}(\hat{Y}_0, Y_0)$  were derived above implies that they can also validly serve, each, as an estimate of the conditional  $\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X})$  given by (32a). This will remain true for all the estimates derived for the MSPE under the other scenarios examined in this paper.*

## 4.2. MSPE in the LM Fitted by OLS with $\mathbf{X}$ Not of Full Column Rank

### 4.2.1. LM Fitted by OLS with $\mathbf{X}$ Column Rank Deficient

Although Assumption  $(\mathcal{A}_3)$  is routinely admitted by most people when handling the LM, one actually meets many concrete instances of data sets where it does not hold. Fortunately, with the formalization by Moore [24], Penrose [25] and, especially, Rao [26] of the notion of *generalized inverse* (short: *g-inverse*) of an arbitrary matrix, it became possible to handle least squares estimation in the LM without having to assume the design matrix  $\mathbf{X}$  necessarily of full column rank.

To begin with, it is shown in most textbooks on the LM that whatever the rank of the design matrix  $\mathbf{X}$ , a vector  $\hat{\beta} \in \mathbb{R}^p$  is a solution to the OLS minimization problem (28) *if, and only if*,  $\hat{\beta}$  is a solution to the so called *normal equations*:

$$\mathbf{X}^T \mathbf{X} \hat{\beta} = \mathbf{X}^T \mathbf{y}. \tag{37}$$

When Assumption  $(\mathcal{A}_3)$  holds, the unique solution to the normal equations is clearly  $\hat{\beta} = \hat{\beta}_{\text{OLS}}$  given by (29). When that's not the case, the square matrix  $\mathbf{X}^T \mathbf{X}$  is singular, hence does not have a regular inverse  $(\mathbf{X}^T \mathbf{X})^{-1}$ . Nevertheless, even then it can be shown that the normal Equation (37) are always consistent. But the apparent negative thing is that they then have infinitely many solution vectors  $\hat{\beta}$ , actually all vectors  $\hat{\beta} \in \mathbb{R}^p$  of the form:

$$\hat{\beta} = \hat{\beta}_{\text{OLS}}^- = (\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{y}, \tag{38}$$

where  $(\mathbf{X}^T \mathbf{X})^-$  is any *g-inverse* of  $\mathbf{X}^T \mathbf{X}$  in the sense of Rao. Given that multitude of possible OLS estimates of  $\beta$  in this case, one may worry that this may hinder any attempt to get a meaningful estimate of the MSPE in the fitted LM. But we are going to show that such a worry is not warranted.

When Assumption  $(\mathcal{A}_3)$  does not hold, in spite of there being as many solutions  $\hat{\beta}_{\text{OLS}}^-$  to the normal Equation (37) as there are *g-inverse*s  $(\mathbf{X}^T \mathbf{X})^-$  of

$\mathbf{X}^T \mathbf{X}$ , *i.e.* infinitely many, it is a remarkable well known fact that the fitted response vector

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\beta}_{OLS}^- = H_X \cdot \mathbf{y}, \quad \text{with } H_X = \mathbf{X}(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T, \quad (39)$$

is the same, whichever *g*-inverse  $(\mathbf{X}^T \mathbf{X})^-$  of  $\mathbf{X}^T \mathbf{X}$  is used to compute  $\hat{\beta}_{OLS}^-$  through (38). This stems from the hat matrix  $H_X$  being equal to the matrix of the orthogonal projection of  $\mathbb{R}^n$  into the range space of  $\mathbf{X}$  ([22] Appendix A). Therefore, the residual response vector

$$\hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X} \hat{\beta}_{OLS}^- = (\mathbf{I}_n - H_X) \mathbf{y}, \quad (40)$$

does not vary with the *g*-inverse  $(\mathbf{X}^T \mathbf{X})^-$  either. Then we will need the result:

**Lemma 4** *With  $r_X = \text{rank}(\mathbf{X})$ , one has:  $\text{tr}(H_X) = r_X$  and  $\mathbb{E}(\|\hat{\mathbf{e}}\|^2 | \mathbf{X}) = (n - r_X) \cdot \sigma^2$ .*

*Proof.* On the one hand, one has:

$$\text{tr}(H_X) = \text{tr}[\mathbf{X}(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T] = \text{tr}[(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{X}]. \quad (41)$$

Now,  $(\mathbf{X}^T \mathbf{X})^-$  being a *g*-inverse of  $\mathbf{X}^T \mathbf{X}$ , it comes that  $(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{X}$  is an idempotent matrix ([22] Appendix A, page 509]. Hence

$$\text{tr}[(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{X}] = \text{rank}[(\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{X}] = \text{rank}(\mathbf{X}^T \mathbf{X}) = \text{rank}(\mathbf{X}) = r_X. \quad (42)$$

Relations (41) and (42) give  $\text{tr}(H_X) = r_X$ . On the other hand,

$$\hat{\mathbf{e}} = (\mathbf{I}_n - H_X) \mathbf{y} = (\mathbf{I}_n - H_X)(\mathbf{X} \beta + \varepsilon) = (\mathbf{I}_n - H_X) \mathbf{X} \beta + (\mathbf{I}_n - H_X) \varepsilon. \quad (43)$$

Now,  $H_X$  being the matrix of the orthogonal projection of  $\mathbb{R}^n$  into the range space of  $\mathbf{X}$ , then  $H_X \mathbf{X} = \mathbf{X}$ . Hence,

$$(\mathbf{I}_n - H_X) \mathbf{X} = \mathbf{X} - H_X \mathbf{X} = 0. \quad (44)$$

From (43) and (44), we get  $\hat{\mathbf{e}} = (\mathbf{I}_n - H_X) \varepsilon$ . Therefore:

$$\|\hat{\mathbf{e}}\|^2 = \varepsilon^T (\mathbf{I}_n - H_X)^T (\mathbf{I}_n - H_X) \varepsilon = \varepsilon^T (\mathbf{I}_n - H_X)^2 \varepsilon = \varepsilon^T (\mathbf{I}_n - H_X) \varepsilon.$$

Under Assumptions  $(\mathcal{A}_0), (\mathcal{A}_1)$ , and thanks to the known identity which gives the expectation of a quadratic form ([27] Appendix B, page 170], one has:

$$\begin{aligned} \mathbb{E}[\|\hat{\mathbf{e}}\|^2 | \mathbf{X}] &= [\mathbb{E}(\varepsilon | \mathbf{X})]^T (\mathbf{I}_n - H_X) \mathbb{E}(\varepsilon | \mathbf{X}) + \text{tr}[(\mathbf{I}_n - H_X) \text{Mcov}(\varepsilon | \mathbf{X})] \\ &= \text{tr}[(\mathbf{I}_n - H_X) \sigma^2 \mathbf{I}_n] = \sigma^2 [\text{tr}(\mathbf{I}_n) - \text{tr}(H_X)] = \sigma^2 (n - r_X) \end{aligned}$$

□

An unbiased estimate of the LM residual variance  $\sigma^2$  in this case is thus known to be:

$$\hat{\sigma}_X^2 = \frac{\text{SSR}}{n - r_X}, \quad \text{with } \text{SSR} = \|\hat{\mathbf{e}}\|^2 = \|\mathbf{y} - \mathbf{X} \hat{\beta}_{OLS}^-\|^2 = \sum_{i=1}^n (y_i - \mathbf{x}_i^T \hat{\beta}_{OLS}^-)^2. \quad (45)$$

We will also need the mean vector and covariance matrix of  $\hat{\beta}_{OLS}^-$ . First,

$$\mathbb{E}(\hat{\beta}_{OLS}^- | \mathbf{X}) = (\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{X} \beta, \quad (46a)$$

which shows that  $\hat{\beta}_{OLS}^-$  is a biased estimator of  $\beta$  when Assumption  $(\mathcal{A}_3)$

does not hold. But in spite of that, note that from (39),

$$\mathbb{E}(\hat{\mathbf{y}} | \mathbf{X}) = \mathbb{E}(\mathbf{X}\hat{\beta}_{OLS}^- | \mathbf{X}) = \mathbf{X}\mathbb{E}(\hat{\beta}_{OLS}^- | \mathbf{X}) = H_X \mathbf{X}\beta = \mathbf{X}\beta = \mathbb{E}(\mathbf{y} | \mathbf{X}). \quad (46b)$$

On the other hand,

$$\text{Mlcov}(\hat{\beta}_{OLS}^- | \mathbf{X}) = \sigma^2 (\mathbf{X}^T \mathbf{X})_S^-, \text{ with } (\mathbf{X}^T \mathbf{X})_S^- = (\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{X} \left[ (\mathbf{X}^T \mathbf{X})^- \right]^T, \quad (47a)$$

the symmetric and positive semi-definite matrix  $(\mathbf{X}^T \mathbf{X})_S^-$  also being a  $g$ -inverse of  $\mathbf{X}^T \mathbf{X}$ . Then

$$\begin{aligned} \text{Mlcov}(\hat{\mathbf{y}} | \mathbf{X}) &= \text{Mlcov}(\mathbf{X}\hat{\beta}_{OLS}^- | \mathbf{X}) = \mathbf{X}\text{Mlcov}(\hat{\beta}_{OLS}^- | \mathbf{X})\mathbf{X}^T \\ &= \sigma^2 \mathbf{X}(\mathbf{X}^T \mathbf{X})_S^- \mathbf{X}^T = \sigma^2 H_X, \end{aligned} \quad (47b)$$

again independent of the  $g$ -inverse  $(\mathbf{X}^T \mathbf{X})^-$  of  $\mathbf{X}^T \mathbf{X}$  used to compute  $\hat{\beta}_{OLS}^-$ .

#### 4.2.2. Preliminary for the MSPE in the LM Fitted by OLS without Assumption ( $\mathcal{A}_3$ )

Our first aim here is to examine the MSPE in the LM when fitted by OLS under the assumption that the design matrix  $\mathbf{X}$  might not have full column rank. So  $\beta$  has been estimated through  $\hat{\beta}_{OLS}^-$  given by (38) and we are interested in  $\text{MSPE}(\hat{Y}_0, Y_0) = \mathbb{E}(\hat{Y}_0 - Y_0)^2$ , where  $\hat{Y}_0 = X_0^T \hat{\beta}_{OLS}^-$  is taken as prediction of  $Y$  on an independently sampled new individual for whom  $X = X_0 \in \mathbb{R}^p$  would have been observed, but not  $Y = Y_0 \in \mathbb{R}$ . We are going to use the results of Section 3.2. First note that since, from the above,  $\hat{\beta}_{OLS}^-$  is a biased estimator of  $\beta$ , Corollary 2 does not apply here. Nonetheless, from Theorem 1, we get:

$$\text{MSPE}(\hat{Y}_0, Y_0 | \mathbf{X}) = \sigma^2 + \text{tr} \left[ \text{Mlse}(\hat{\beta}_{OLS}^- | \mathbf{X}) \cdot \mathbb{E}(X_0 X_0^T) \right], \quad (48a)$$

$$\text{MSPE}(\hat{Y}_0, Y_0) = \sigma^2 + \text{tr} \left[ \text{Mlse}(\hat{\beta}_{OLS}^-) \cdot \mathbb{E}(X_0 X_0^T) \right]. \quad (48b)$$

Our estimation of  $\text{MSPE}(\hat{Y}_0, Y_0)$  in this case will be based on those two identities and:

**Lemma 5** *In the LM(14) with Assumptions ( $\mathcal{A}_0$ )-( $\mathcal{A}_1$ ),*

$$\text{tr} \left[ \text{Mlse}(\hat{\beta}_{OLS}^- | \mathbf{X}) \cdot \mathbf{X}^T \mathbf{X} \right] = r_X \cdot \sigma^2, \text{ with } r_X = \text{rank}(\mathbf{X}). \quad (49)$$

*Proof.* For

$$A(\mathbf{X}) = \text{tr} \left[ \text{Mlse}(\hat{\beta}_{OLS}^- | \mathbf{X}) \cdot \mathbf{X}^T \mathbf{X} \right] = \text{tr} \left\{ \mathbb{E} \left[ (\hat{\beta}_{OLS}^- - \beta)(\hat{\beta}_{OLS}^- - \beta)^T | \mathbf{X} \right] \cdot \mathbf{X}^T \mathbf{X} \right\},$$

$$A(\mathbf{X}) = \text{tr} \left\{ \mathbf{X} \cdot \mathbb{E} \left[ (\hat{\beta}_{OLS}^- - \beta)(\hat{\beta}_{OLS}^- - \beta)^T | \mathbf{X} \right] \cdot \mathbf{X}^T \right\}$$

$$= \text{tr} \left\{ \mathbb{E} \left[ \mathbf{X}(\hat{\beta}_{OLS}^- - \beta)(\hat{\beta}_{OLS}^- - \beta)^T \mathbf{X}^T | \mathbf{X} \right] \right\}$$

$$= \text{tr} \left\{ \mathbb{E} \left[ (\hat{\mathbf{y}} - \mathbf{X}\beta)(\hat{\mathbf{y}} - \mathbf{X}\beta)^T | \mathbf{X} \right] \right\}$$

$$= \text{tr} \left[ \text{Mlcov}(\hat{\mathbf{y}} | \mathbf{X}) \right], \text{ thanks to (46b)}$$

$$= \sigma^2 \text{tr}(H_X), \text{ by (47b).}$$

Hence (49), thanks to Lemma 4.  $\square$



### 4.2.3. The FPE and the GCV in the LM Fitted by OLS with $X$ Column Rank Deficient

Given (48a), we first estimate  $\mathbb{E}(X_0 X_0^T)$  by  $X^T X/n$  as in (33), which entails the following preliminary estimate of  $\text{MSPE}(\hat{Y}_0, Y_0 | X)$  in the present case, using the last lemma:

$$\widehat{\text{MSPE}}(\hat{Y}_0, Y_0 | X) = \sigma^2 + \frac{1}{n} \text{tr} \left[ \text{Mlse}(\hat{\beta}_{\text{OLS}} | X) \cdot X^T X \right] = \sigma^2 \left( 1 + \frac{r_X}{n} \right). \quad (50)$$

But since  $\text{MSPE}(\hat{Y}_0, Y_0) = \mathbb{E}_X \text{MSPE}(\hat{Y}_0, Y_0 | X)$ , then (51) also gives a preliminary estimate of  $\text{MSPE}(\hat{Y}_0, Y_0)$ . Then, estimating  $\sigma^2$  by  $\hat{\sigma}_X^2$  given by (45), our final estimate of  $\text{MSPE}(\hat{Y}_0, Y_0)$  in this case, computable from data, is:

$$\text{FPE} = \hat{\sigma}_X^2 \left( 1 + \frac{r_X}{n} \right) = \hat{S}^2 \cdot \frac{n + r_X}{n - r_X}, \quad (51)$$

which is also an estimate of  $\text{MSPE}(\hat{Y}_0, Y_0 | X)$ . It is denoted FPE because it generalizes (34) in assessing the goodness-of-fit of OLS in the LM when the design matrix  $X$  is column rank deficient. The remarkable feature is that this estimate is the same, whichever  $g$ -inverse  $(X^T X)^-$  of  $X^T X$  was used to get the estimate  $\hat{\beta}_{\text{OLS}}$  of  $\beta$  in (38).

Estimating, instead,  $\mathbb{E}(X_0 X_0^T)$  by  $X^T X / (n - r_X)$  yields as estimate of  $\text{MSPE}(\hat{Y}_0, Y_0)$ :

$$\hat{\sigma}_X^2 \left( 1 + \frac{r_X}{n - r_X} \right) = \frac{n \cdot \text{SSR}}{(n - r_X)^2} = \frac{\hat{S}^2}{(1 - r_X/n)^2} = \text{GCV}, \quad (52)$$

the GCV estimate of  $\text{MSPE}(\hat{Y}_0, Y_0)$  in the LM fitted by OLS when  $X$  is not full column rank.

## 5. MSPE when Fitting the LM by Ridge Regression

The design matrix  $X$  being column rank deficient means that its  $p$  columns are linearly dependent, or almost so. This happens when there is multicollinearity among the  $p$  regressors  $X_1, \dots, X_p$ , and thus some of them are redundant with some others. However, when this occurs, computing an OLS estimate  $\hat{\beta}_{\text{OLS}}$  of  $\beta$ , given by (38), in a numerically stable manner is not easy and requires using carefully designed Numerical Linear Algebra programming routines. The difficulty stems from the fact that this requires, at least implicitly, to find, along the way, the exact rank of  $X$ , which is difficult to achieve, precisely because of the multicollinearity among its columns. It can then be of much interest to have a method which can fit the LM without having to bother about the exact rank of  $X$ . This is precisely what Ridge regression (RR) tries to achieve.

Hoerl and Kennard presented two variants of Ridge Regression [19]. In the initial one (the default), which some have termed Ordinary Ridge Regression (ORR), to fit the LM (14), one estimates  $\beta$  through regularizing the OLS criterion (28) by a ridge constraint, yielding:

$$\hat{\beta} = \hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \left[ \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 + \lambda \|\beta\|^2 \right], \quad (53)$$

for some  $\lambda > 0$ , a penalty parameter to choose appropriately. The unique solution to (54) is known to be (whatever the rank of  $\mathbf{X}$ ):

$$\hat{\beta}_\lambda = G_\lambda^{-1} \mathbf{X}^T \mathbf{y}, \quad \text{with} \quad G_\lambda = \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p. \tag{54}$$

Hoerl and Kennard presented ORR in [19] assuming that the design matrix  $\mathbf{X}$  was full column rank (*i.e.* our Assumption ( $\mathcal{A}_3$ )), which also requires that  $n \geq p$ . But because the symmetric matrix  $\mathbf{X}^T \mathbf{X}$  is always at least semi-positive definite, imposing  $\lambda > 0$  entails that the  $p \times p$  matrix  $G_\lambda$  is symmetric and positive definite (SPD), whatever the rank of  $\mathbf{X}$  and the ranking between the integers  $n$  and  $p$ . That is why, in what follows, we do not impose any rank constraint on  $\mathbf{X}$  (apart being trivially  $\geq 1$  because  $\mathbf{X}$  is nonzero). No specific ranking either is assumed between  $n$  and  $p$ .

However, hereafter, since it does not require extra work, we directly consider the extended setting of Generalized Ridge Regression (GRR) which, to fit the LM (14), estimates  $\beta$  through solving the minimization problem:

$$\hat{\beta} = \hat{\beta}_\lambda = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \left[ \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 + \lambda \cdot \|\beta\|_D^2 \right], \tag{55}$$

where  $\lambda$  is as in ORR and  $\mathbf{D}$  is a  $p \times p$  symmetric and semi-positive definite (SSPD) matrix, both given, while  $\|\beta\|_D^2 = \beta^T \mathbf{D} \beta$ . The solution of (56) is still (55), but now with

$$G_\lambda = \mathbf{X}^T \mathbf{X} + \lambda \mathbf{D}, \tag{56}$$

under the assumption that the latter is SPD. Since smoothing splines fitting can be cast in the GRR form (56), what follows applies to that hugely popular non-parametric regression method as well.

### 5.1. The MSPE Issue for Ridge Regression

More than for Least Squares, depending on the unspecified parameter  $\lambda$ , it is critical to assess how Ridge Regression fits the LM for given data in order to be able to select the best  $\lambda$  value, *i.e.* the one ensuring the best fit. From the prediction error point of view stated in this article, this amounts to choosing the  $\lambda > 0$  for which the RR fit has the smallest MSPE. It, thus, requires to estimate  $\text{MSPE}_\lambda(\hat{Y}_0, Y_0) = \mathbb{E}(\hat{Y}_0 - Y_0)^2$  for any given  $\lambda$  value, where  $X_0$  and  $Y_0$  are as before, while  $\hat{Y}_0 = X_0^T \hat{\beta}_\lambda$ . Traditionally, estimating  $\text{MSPE}_\lambda(\hat{Y}_0, Y_0)$  in this context is mostly done using the Generalized Cross Validation (GCV) criterion initially developed by Craven and Wahba for selecting the best value of the smoothing parameter in a smoothing spline [3]. That GCV is obtained as a variation of the LOO-CV. Here, to estimate  $\text{MSPE}_\lambda(\hat{Y}_0, Y_0)$ , we take a different route. We first note that the fitted sample response vector is

$$\hat{\mathbf{y}}_\lambda = \mathbf{X} \hat{\beta}_\lambda = H_\lambda \mathbf{y}, \quad \text{with} \quad H_\lambda = \mathbf{X} G_\lambda^{-1} \mathbf{X}^T. \tag{57}$$

On the other hand, from (55),

$$\mathbb{E}(\hat{\beta}_\lambda | \mathbf{X}) = G_\lambda^{-1} \mathbf{X}^T \mathbb{E}(\mathbf{y} | \mathbf{X}) = T_\lambda \beta, \quad \text{with} \quad T_\lambda = G_\lambda^{-1} \mathbf{X}^T \mathbf{X} \neq \mathbf{I}_p. \tag{58}$$

So, again, Corollary 2 does not apply. But, using Theorem 1,

$$\text{MSPE}_\lambda(\hat{Y}_0, Y_0 | \mathbf{X}) = \sigma^2 + \text{tr} \left[ \text{Mlse}(\hat{\beta}_\lambda | \mathbf{X}) \cdot \mathbb{E}(X_0 X_0^T) \right], \tag{59a}$$

$$\text{MSPE}_\lambda(\hat{Y}_0, Y_0) = \sigma^2 + \text{tr} \left[ \text{Mlse}(\hat{\beta}_\lambda) \cdot \mathbb{E}(X_0 X_0^T) \right]. \tag{59b}$$

For estimating those quantities, we will need some preliminary results.

### 5.2. Preliminary Results for Estimating the MSPE in Ridge Regression

First, two simple, but remarkable identities about the matrices  $H_\lambda$  and  $T_\lambda$  given in (58) and (59).

**Lemma 6**  $T_\lambda = \mathbf{I}_p - \lambda G_\lambda^{-1} \mathbf{D}$  and  $\lambda \mathbf{X} G_\lambda^{-1} \mathbf{D} = (\mathbf{I}_n - H_\lambda) \mathbf{X}$ .

*Proof.* The first identity is easily got from (59) and (57). Indeed, one has:

$$\mathbf{I}_p = G_\lambda^{-1} G_\lambda = G_\lambda^{-1} \mathbf{X}^T \mathbf{X} + G_\lambda^{-1} \lambda \mathbf{D} = T_\lambda + \lambda G_\lambda^{-1} \mathbf{D},$$

As for the second one,

$$\lambda \mathbf{X} G_\lambda^{-1} \mathbf{D} = \mathbf{X} G_\lambda^{-1} (G_\lambda - \mathbf{X}^T \mathbf{X}) = \mathbf{X} - \mathbf{X} G_\lambda^{-1} \mathbf{X}^T \mathbf{X} = (\mathbf{I}_n - H_\lambda) \mathbf{X}. \quad \square$$

□

Next, a key preliminary about the Mean Squared Error matrix of  $\hat{\beta}_\lambda$  as an estimate of  $\beta$ :

**Lemma 7** Under Assumptions  $(\mathcal{A}_0)$ - $(\mathcal{A}_1)$ , one has:

$$\text{tr} \left[ \text{Mlse}(\hat{\beta}_\lambda | \mathbf{X}) \cdot \mathbf{X}^T \mathbf{X} \right] = \sigma^2 \text{tr}(H_\lambda^2) + \|(\mathbf{I}_n - H_\lambda) \mathbf{X} \beta\|^2. \tag{60}$$

*Proof.* Inserting (20) in  $\text{tr} \left[ \text{Mlse}(\hat{\beta}_\lambda | \mathbf{X}) \cdot \mathbf{X}^T \mathbf{X} \right] = \text{tr} \left[ \mathbf{X} \text{Mlse}(\hat{\beta}_\lambda | \mathbf{X}) \mathbf{X}^T \right]$  gives:

$$\text{tr} \left[ \text{Mlse}(\hat{\beta}_\lambda | \mathbf{X}) \cdot \mathbf{X}^T \mathbf{X} \right] = \text{tr}(A) + \text{tr}(B),$$

with  $A = \mathbf{X} \text{Mlcov}(\hat{\beta}_\lambda | \mathbf{X}) \mathbf{X}^T$  and  $B = \mathbf{X} \text{Bias}(\hat{\beta}_\lambda) \cdot \left[ \mathbf{X} \text{Bias}(\hat{\beta}_\lambda) \right]^T$ . Now,

$$A = \text{Mlcov}(\hat{\mathbf{y}}_\lambda | \mathbf{X}) = H_\lambda \text{Mlcov}(\mathbf{y} | \mathbf{X}) H_\lambda = \sigma^2 H_\lambda^2$$

$$\begin{aligned} \text{tr}(B) &= \text{tr} \left\{ \left[ \mathbf{X} \text{Bias}(\hat{\beta}_\lambda) \right]^T \mathbf{X} \text{Bias}(\hat{\beta}_\lambda) \right\} = \left\| \mathbf{X} \text{Bias}(\hat{\beta}_\lambda) \right\|^2 \\ &= \left\| \mathbf{X} (T_\lambda - \mathbf{I}_p) \beta \right\|^2 = \left\| \lambda \mathbf{X} G_\lambda^{-1} \mathbf{D} \beta \right\|^2 = \left\| (\mathbf{I}_n - H_\lambda) \mathbf{X} \beta \right\|^2, \end{aligned}$$

the last three identities using (59) and Lemma 6. □

With the above lemma, we are now in a position to be able to estimate  $\text{MSPE}_\lambda(\hat{Y}_0, Y_0)$  in Ridge Regression. We examine, hereafter, two paths for achieving that: one leads to the FPE, the other one to the GCV.

### 5.3. Estimating the MSPE in Ridge Regression by the FPE

Here, we first estimate  $\mathbb{E}(X_0 X_0^T)$  by  $\mathbf{X}^T \mathbf{X} / n$  in (60a). Then, given (61), this suggests the preliminary estimate of  $\text{MSPE}_\lambda(\hat{Y}_0, Y_0 | \mathbf{X})$  in RR:

$$\begin{aligned} \widehat{\text{MSPE}}_{\lambda}^0(\hat{Y}_0, Y_0 | \mathbf{X}) &= \sigma^2 + \frac{1}{n} \text{tr} \left[ \text{Mise}(\hat{\beta}_{\lambda} | \mathbf{X}) \cdot \mathbf{X}^T \mathbf{X} \right] \\ &= \sigma^2 \left[ 1 + \frac{\text{tr}(H_{\lambda}^2)}{n} \right] + \frac{\|(\mathbf{I}_n - H_{\lambda}) \mathbf{X} \beta\|^2}{n}, \end{aligned} \tag{61}$$

which is also, therefore, a preliminary estimate of  $\text{MSPE}_{\lambda}(\hat{Y}_0, Y_0)$ . It is only a preliminary estimate because even given  $\lambda$ , it still depends on the two unknowns  $\sigma^2$  and  $\beta$ . Interestingly, it can be shown ([8] Chapter 3, page 46) that (62) coincides with  $\text{PSE}(\hat{y}_{\lambda} | \mathbf{X})$  given by (62) in the present setting.

It is useful to note that (62) depends on  $\beta$  only through the squared bias term  $\|(\mathbf{I}_n - H_{\lambda}) \mathbf{X} \beta\|^2$ . To estimate the latter, let  $\hat{e}_{\lambda} = \mathbf{y} - \hat{y}_{\lambda} = (\mathbf{I}_n - H_{\lambda}) \mathbf{y}$ , the vector of sample residuals, and  $\text{SSR}_{\lambda} = \|\hat{e}_{\lambda}\|^2$ , the sum of squared residuals in the Ridge Regression fit. Then, thanks to a well known identity ([9] Chapter 2, page 38),

$$\mathbb{E}(\text{SSR}_{\lambda} | \mathbf{X}) = \|(\mathbf{I}_n - H_{\lambda}) \mathbf{X} \beta\|^2 + \sigma^2 \text{tr} \left[ (\mathbf{I}_n - H_{\lambda})^2 \right],$$

implying that  $\text{SSR}_{\lambda} - \sigma^2 \text{tr} \left[ (\mathbf{I}_n - H_{\lambda})^2 \right]$  is an unbiased estimate of  $\|(\mathbf{I}_n - H_{\lambda}) \mathbf{X} \beta\|^2$ . Hence a general formula for computing an estimate of the MSPE in Ridge Regression:

$$\text{FPE}(\hat{\sigma}^2) = \frac{2\hat{\sigma}^2 \text{tr}(H_{\lambda}) + \text{SSR}_{\lambda}}{n}, \tag{62}$$

where  $\hat{\sigma}^2$  is a chosen estimate of  $\sigma^2$ , possibly computed from the RR fit, thus dependent on  $\lambda$ .

We denoted  $\text{FPE}(\hat{\sigma}^2)$  the estimate of the MSPE given by (63) for the reason to follow. Indeed, probably the most popular estimate of the residual variance  $\sigma^2$  from an RR fit is the one proposed by Wahba in the context of smoothing splines [28]:

$$\hat{\sigma}_1^2 = \frac{\text{SSR}_{\lambda}}{n - \text{tr}(H_{\lambda})}. \tag{63}$$

Now, if one uses  $\hat{\sigma}^2 = \hat{\sigma}_1^2$  in (63), an algebraic manipulation easily leads to:

$$\text{FPE}(\hat{\sigma}_1^2) = \frac{\text{SSR}_{\lambda}}{n} \cdot \frac{n + \text{tr}(H_{\lambda})}{n - \text{tr}(H_{\lambda})} = \text{FPE}, \tag{64}$$

recovering the classical formula of the FPE for this setting, but this time as an estimate of the MSPE rather than the PSE.

### 5.4. Estimating the MSPE in Ridge Regression by the GCV

Here, we estimate  $\mathbb{E}(X_0 X_0^T)$  in (60a) rather by  $\mathbf{X}^T \mathbf{X} / [n - \text{tr}(H_{\lambda})]$ . With (61), this suggests a preliminary estimate  $\widehat{\text{MSPE}}_{\lambda}^1(\hat{Y}_0, Y_0 | \mathbf{X})$  of  $\text{MSPE}_{\lambda}(\hat{Y}_0, Y_0 | \mathbf{X})$  in RR, corresponding to  $\widehat{\text{MSPE}}_{\lambda}^0(\hat{Y}_0, Y_0 | \mathbf{X})$  where the denominators  $n$  have been replaced by  $n - \text{tr}(H_{\lambda})$ . Using the same unbiased estimate of  $\|(\mathbf{I}_n - H_{\lambda}) \mathbf{X} \beta\|^2$  as in the previous section and an estimate  $\hat{\sigma}^2$  of  $\sigma^2$ , we get another general formula for computing an estimate of the MSPE in RR:

$$\text{GCV}(\hat{\sigma}^2) = \frac{\hat{\sigma}^2 \text{tr}(H_\lambda) + \text{SSR}_\lambda}{n - \text{tr}(H_\lambda)}. \quad (65)$$

We denoted it  $\text{GCV}(\hat{\sigma}^2)$  because again taking  $\hat{\sigma}^2 = \hat{\sigma}_1^2$ , one gets:

$$\text{GCV}(\hat{\sigma}_1^2) = \frac{\text{SSR}_\lambda}{n[1 - n^{-1}\text{tr}(H_\lambda)]^2} = \text{GCV}, \quad (66)$$

the well known formula of Craven and Wahba's GCV [3] for this setting, but this time derived without any reference to Cross Validation.

## 6. Conclusion and Perspectives

In this work, the goal was not to derive new and better selection criteria to assess the goodness-of-fit in regression, but rather to show how one can derive the well known Akaike's FPE and, Craven and Wahba's GCV [3] as direct estimates of the measure of prediction error most commonly known to users, which is not how they are traditionally derived. We achieved this for some of the best known linear model fitters, with the two derivations differing only slightly for each of them. But, nowadays, in regression, use of the FPE criterion is generally not recommended because much better performing criteria are known [29], while GCV has its own shortcomings in certain settings (*e.g.* small sample size), though hugely popular and almost the best in some difficult high dimensional situations [12]. It is then our hope that, in the future, one can, through the same unified framework used in this paper, derive new and better selection criteria, different from other already available such proposals for the same setting, among which we can cite the AICc [30], the modified GCV [31] [32], the modified RGCV and  $R_1\text{GCV}$  [33] [34],  $\text{RC}_p$ ,  $\text{RC}_p^+$  and  $\widehat{\text{RC}}_p$  [5].

## Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

## References

- [1] Akaike, H. (1970) Statistical Predictor Identification. *Annals of the Institute of Statistical Mathematics*, **22**, 203-217. <https://doi.org/10.1007/BF02506337>
- [2] Borra, S. and Di Ciaccio, A. (2010) Measuring the Prediction Error. A Comparison of Cross-Validation, Bootstrap and Covariance Penalty Methods. *Computational Statistics & Data Analysis*, **54**, 2976-2989. <https://doi.org/10.1016/j.csda.2010.03.004>
- [3] Craven, P. and Wahba, G. (1979) Smoothing Noisy Data with Spline Functions: Estimating the Correct Degree of Smoothing by the Method of Generalized Cross-Validation. *Numerische Mathematik*, **31**, 377-403. <https://doi.org/10.1007/BF01404567>
- [4] Li, K.-C. (1985) From Stein's Unbiased Risk Estimates to the Method of Generalized Cross Validation. *Journal of the Japan Statistical Society*, **38**, 119-130.
- [5] Rosset, S. and Tibshirani, R. (2020) From Fixed-X to Random-X Regression: Bias-Variance Decompositions, Covariance Penalties, and Prediction Error Estima-

- tion. *JASA*, **115**, 138-151. <https://doi.org/10.1080/01621459.2018.1424632>
- [6] Hastie, T., Tibshirani, R. and Friedman, J. (2009) *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. 2nd Edition, Springer-Verlag, New York.
- [7] Burman, P. (1989) A Comparative Study of Ordinary Cross-Validation, v-Fold Cross-Validation and Repeated Learning-Testing Methods. *Biometrika*, **76**, 503-514. <https://doi.org/10.1093/biomet/76.3.503>
- [8] Hastie, T. and Tibshirani, R. (1990) *Generalized Additive Models*. Monographs on Statistics and Applied Probability. Chapman and Hall, London.
- [9] Eubank, R.L. (1999) *Nonparametric Regression and Spline Smoothing*. 2nd Edition, Marcel Dekker, New York. <https://doi.org/10.1201/9781482273144>
- [10] McQuarrie, A.D.R. and Tsai, C.-L. (1998) *Regression and Time Series Model Selection*. World Scientific Publishing Co. Re. Ltd, Singapore. <https://doi.org/10.1142/3573>
- [11] Breiman, L. and Spector, P. (1992) Submodel Selection and Evaluation in Regression. The X-Random Case. *International Statistical Review*, **60**, 291-319. <https://doi.org/10.2307/1403680>
- [12] Leeb, H. (2008) Evaluation and Selection of Models for Out-of-Sample Prediction When the Sample Size Is Small Relative to the Complexity of the Data-Generating Process. *Bernoulli*, **14**, 661-690. <https://doi.org/10.3150/08-BEJ127>
- [13] Dicker, L.H. (2013) Optimal Equivariant Prediction for High-Dimensional Linear Models with Arbitrary Predictor Covariance. *The Electronic Journal of Statistics*, **7**, 1806-1834. <https://doi.org/10.1214/13-EJS826>
- [14] Dobriban, E. and Wager, S. (2018) High-Dimensional Asymptotics of Prediction. *Annals of Statistics*, **46**, 247-279. <https://doi.org/10.1214/17-AOS1549>
- [15] Lukas, M.A. (2014) Performance Criteria and Discrimination of Extreme Undersmoothing in Nonparametric Regression. *Journal of Statistical Planning and Inference*, **153**, 56-74. <https://doi.org/10.1016/j.jspi.2014.05.006>
- [16] Lukas, M.A., de Hoog, F.R. and Anderssen, R.S. (2015) Practical Use of Robust GCV and Modified GCV for Spline Smoothing. *Computational Statistics*, **31**, 269-289. <https://doi.org/10.1007/s00180-015-0577-7>
- [17] Kitagawa, G. (2008) Contributions of Professor Hirotogu Akaike in Statistical Science. *Journal of the Japan Statistical Society*, **38**, 119-130. <https://doi.org/10.14490/jjss.38.119>
- [18] Tibshirani, R. (1996) Regression Shrinkage and Selection via the Lasso. *JRSS, Series B*, **58**, 267-288. <https://doi.org/10.1111/j.2517-6161.1996.tb02080.x>
- [19] Hoerl, A.E. and Kennard, R. (1970) Ridge Regression: Biased Estimation for Non-orthogonal Problems. *Technometrics*, **12**, 55-67. <https://doi.org/10.1080/00401706.1970.10488634>
- [20] Searle, S.R. (1971) *Linear Models*. John Wiley & Sons, New York.
- [21] Rao, C.R. (1973) *Linear Statistical Inference and Its Applications*. 2nd Edition, Wiley, New York. <https://doi.org/10.1002/9780470316436>
- [22] Rao, C.R., Toutenburg, H., *et al.* (2008) *Linear Models and Generalizations: Least Squares and Alternatives*. 3rd Edition, Springer-Verlag, Berlin.
- [23] Olive, D.J. (2017) *Linear Regression*. Springer International Publishing, Berlin. <https://doi.org/10.1007/978-3-319-55252-1>
- [24] Moore, E.H. (1920) On the Reciprocal of the General Algebraic Matrix (Abstract).

*Bulletin of the AMS*, **26**, 394-395.

- [25] Penrose, R. (1955) A Generalized Inverse for Matrices. *Proceedings—Cambridge Philosophical Society*, **51**, 406-413. <https://doi.org/10.1017/S0305004100030401>
- [26] Rao, C.R. (1962) A Note on a Generalized Inverse of a Matrix with Applications to Problems in Mathematical Statistics. *JRSS, Series B*, **24**, 152-158. <https://doi.org/10.1111/j.2517-6161.1962.tb00447.x>
- [27] Kendrick, D.A. (2002) Stochastic Control for Economic Models. 2nd Edition, VTEX Ltd., Vilnius.
- [28] Wahba, G. (1978) Improper Priors, Spline Smoothing and the Problem of Guarding against Model Errors in Regression. *JRSS, Series B*, **40**, 364-372. <https://doi.org/10.1111/j.2517-6161.1978.tb01050.x>
- [29] Härdle, W., Hall, P. and Marron, J.S. (1988) How Far Are Automatically Chosen Regression Smoothing Parameters from Their Optimum-(with Discussion). *JASA*, **83**, 86-89. <https://doi.org/10.2307/2288922>
- [30] Hurvich, C.M., Simonoff, J.S. and Tsai, C.-L. (1998) Smoothing Parameter Selection in Nonparametric Regression Using an Improved Akaike Information Criterion. *JRSS, Series B*, **60**, 271-293. <https://doi.org/10.1111/1467-9868.00125>
- [31] Cummins, D.J., Filloon, T.G. and Nychka, D. (2001) Confidence Intervals for Non-parametric Curve Estimates: Toward More Uniform Pointwise Coverage. *JASA*, **96**, 233-246. <https://doi.org/10.1198/016214501750332811>
- [32] Kim, Y.-J. and Gu, C. (2004) Smoothing Spline Gaussian Regression: More Scalable Computation via Efficient Approximation. *JRSS, Series B*, **66**, 337-356. <https://doi.org/10.1046/j.1369-7412.2003.05316.x>
- [33] Lukas, M.A. (2006) Robust Generalized Cross-Validation for Choosing the Regularization Parameter. *Inverse Probability*, **22**, 1883-1902. <https://doi.org/10.1088/0266-5611/22/5/021>
- [34] Lukas, M.A. (2008) Strong Robust Generalized Cross-Validation for Choosing the Regularization Parameter. *Inverse Probability*, **24**, Article ID: 034006. <https://doi.org/10.1088/0266-5611/24/3/034006>