# Local Influence for the Thin-Plate Spline Generalized Linear Model

Germán Ibacache-Pulgar [1,2,]*[ID], Pablo Pacheco [3][ID], Orietta Nicolis [4][ID] and Miguel Angel Uribe-Opazo [5][ID]

1 Institute of Statistics, Universidad de Valparaíso, Av. Gran Bretaña 1111, Valparaíso 2360102, Chile
2 Centro de Estudios Atmosféricos y Cambio Climático (CEACC), Universidad de Valparaíso, Valparaíso 2360102, Chile
3 Dirección de Educación Virtual, Universidad de Playa Ancha, Avenida Guillermo González de Hontaneda 855, Playa Ancha, Valparaíso 2360072, Chile; ppachecog@upla.cl
4 Facultad de Ingenieria, Universidad Andres Bello, Calle Quillota 980, Viña del Mar 2520000, Chile; orietta.nicolis@unab.cl
5 Centro de Ciências Exatas e Tecnológicas, Western Paraná State University (UNIOESTE), Cascavel 85819-110, Paraná, Brazil; miguel.opazo@unioeste.br
* Correspondence: german.ibacache@uv.cl

**Abstract:** Thin-Plate Spline Generalized Linear Models (TPS-GLMs) are an extension of Semiparametric Generalized Linear Models (SGLMs), because they allow a smoothing spline to be extended to two or more dimensions. This class of models allows modeling a set of data in which it is desired to incorporate the non-linear joint effects of some covariates to explain the variability of a certain variable of interest. In the spatial context, these models are quite useful, since they allow the effects of locations to be included, both in trend and dispersion, using a smooth surface. In this work, we extend the local influence technique for the TPS-GLM model in order to evaluate the sensitivity of the maximum penalized likelihood estimators against small perturbations in the model and data. We fit our model through a joint iterative process based on Fisher Scoring and weighted backfitting algorithms. In addition, we obtained the normal curvature for the case-weight perturbation and response variable additive perturbation schemes, in order to detect influential observations on the model fit. Finally, two data sets from different areas (agronomy and environment) were used to illustrate the methodology proposed here.

**Keywords:** exponential family; smoothing spline; penalized likelihood function; weighted backfitting algorithm; diagnostics measures

**MSC:** 62P12; 62J20; 62G05

## 1. Introduction

Thin-Plate Spline Generalized Linear Models (TPS-GLMs) represent an extension of semiparametric generalized linear models (SGLMs) by enabling the application of smoothing splines in multiple dimensions. These models have the same characteristics of the generalized linear model (GLM), as described by McCullagh and Nelder [1]. Like GLMs, TPS-GLMs can assume a variety of distribution families for the response variable. They also allow for a non-linear relationship between the response variable's mean and the linear predictor via a link function, and they account for non constant variance in the data. Furthermore, the TPS-GLM allow modeling non-linear joint interaction effects due to some covariates, as well as the effects of coordinates in spatial data, making them a useful tool to model dynamic pattern in different scientific areas, such as environment, agronomy, ecology, and so on. Some of the main works related to thin-plate spline technique are Duchon [2,3], Bookstein [4], and Chen et al. [5], while in the context of statistical modeling, Wahba [6], Green and Silverman [7], Wood [8], and Moraga et al. [9], can be mentioned, among others.

However, it is well known that diagnostic analysis is a fundamental process in all statistical modeling for any data set. This analysis allows us to validate the assumptions established about the model in question and identify discrepant observations, and eventually influential ones on the fit of the model. One of the main diagnostic techniques used in GLM and SGLM is local influence. In general, the idea of the local influence technique introduced by Cook [10] is to evaluate the sensitivity of the MLEs when small perturbations are introduced in the assumptions of the model or in the data, both in the response variable and in the explanatory variables. This technique has the advantage, regarding the case elimination technique, that it is not necessary to calculate the estimates of the parameters for each case excluded. In our case, we are interested in developing the local influence technique in the TPS-GLM, in order to detect observations that may have a disproportionate influence on the estimators of both the parametric (regression coefficient) and non-parametric (surface) part of the linear predictor. Such influence may be due, for example, to the fact that each experimental unit contributes differently to the model or that our variable of interest is exposed to a certain modification. In the context of GLM and SGLM, there is empirical evidence that the maximum likelihood estimators (MLEs) and maximum penalized likelihood estimators (PMLEs) are sensitive to this type of situation, and therefore we believe that this sensitivity is also present in the estimators of the TPS-GLM, in particular, in the surface estimator.

Various studies have expanded upon the technique of local influence within different parametric models. Thomas and Cook [11] applied Cook's method of local influence [10] to generalized linear models to assess the impact of minor data perturbations. Ouwens and Beger [12] obtained the normal curvature under a generalized linear model in order to identify influential subjects and/or individual observations. Zhu and Lee [13] developed the local influence technique for incomplete data, and extended such results to generalized linear mixed models (see also Zhu and Lee [14] for further details). Espinheira et al. [15] extended the local influence analysis to beta regression models considering various perturbation scenarios. Rocha and Simas [16] and Ferrari et al. [17] derived the normal curvature considering a beta regression model whose dispersion parameter varies according to the effect of some covariates. Ferreira and Paula [18] developed the local influence approach to partially linear Skew Normal models under different perturbation schemes, and Emami [19] evaluated the sensitivity of Liu penalized least squares estimators using local influence analysis. Most recently, Liu et al. [20] have reported the implementation of influence diagnostics in AR time series models with Skew Normal (SK) distributions.

Within a semiparametric framework, Thomas [21] developed diagnostics for local influence to assess the sensitivity of estimates for the smoothing parameter, which were determined using the cross-validation criterion. Zhu and Lee [14] and Ibacache-Pulgar and Paula [22] introduced measures of local influence to analyze the sensitivity of maximum penalized likelihood estimates in normal and partially linear Student-t models, respectively. Ibacache-Pulgar et al. [23,24] explored local influence curvature within elliptical semiparametric mixed models and symmetric semiparametric additive models. Subsequently, ref. [25] and Ibacache-Pulgar and Reyes [26] further extended local influence measures to normal and elliptical partially varying-coefficient models, respectively. Ibacache-Pulgar et al. [27] developed the local influence method within the context of semiparametric additive beta regression models. Meanwhile, Cavieres et al. [28] calculated the normal curvature to assess the sensitivity of estimators in a thin-plate spline model that incorporates skew normal random errors. Jeldes et al. [29] applied the partially coefficient-varying model with symmetric random errors to air pollution data from the cities of Santiago, Chile, and Lima, Peru. In this context, they carried out an application of the local influence technique to detect influential observations in the model fit. Saavedra-Nievas et al. [30] extended the local influence technique for the spatio-temporal linear model under normal distribution and with separable covariance. Recently, Sánchez et al. [31] obtained the normal curvature for the varying-coefficient quantile regression model under log-symmetric distributions,

and presented an interesting application of such results to an environmental pollution data set.

In this work, we extend the local influence approach in Thin-Plate Spline Generalized Linear Model.

The contents are organized as follows: Section 2 introduces the thin-plate spline generalized linear model. Section 3 details the method for obtaining maximum penalized likelihood estimators and discusses some statistical inferential results. In Section 4, we provide a detailed description of the local influence method and derives normal curvatures for various perturbation schemes. In Section 5, the methodology is illustrated using two datasets. The paper concludes with some final observations in Section 6.

## 2. The Thin-Plate Spline Generalized Linear Model (TPS-GLM)

In this section, we present the TPS-GLM and the penalized function to carry out the process of estimating the parameters.

### 2.1. Statistical Model

Let $\{y_i \mid i = 1, \ldots, n\}$ be a data set where each response variable $y_i$ follows a distribution from the exponential family with the following density function:

$$f_y(y_i; \theta_i, \phi) = \exp \left[ \frac{y_i \theta_i - \psi(\theta_i)}{a_i(\phi)} + c(y_i, \phi) \right],$$

where $\theta_i$ is the canonical form of the location parameter and depends on the mean $\mu_i$. The term $a_i(\phi)$ represents a known function of the unknown dispersion parameter $\phi$ (or a vector of unknown dispersion parameters). The function $c$ depends on both the dispersion parameter and the responses, while $\psi$ is a known function, such that the mean and variance of $y_i$ are given by: $\mu_i = \mathrm{E}(y_i) = \partial \psi(\theta_i)/\partial \theta_i$ and $\mathrm{Var}(y_i) = a_i(\phi) \mathrm{V}_i$, with $\mathrm{V}_i = \mathrm{V}(\mu_i) = \partial^2 \psi(\theta_i)/\partial \theta_i^2$, respectively. The TPS-GLM is defined by Equation (1) and the following systematic component:

$$g(\mu_i) = \eta_i = \mathbf{w}_i^\top \boldsymbol{\alpha} + f(\mathbf{t}_i), \tag{1}$$

where $\mathbf{w}_i$ is a $(p \times 1)$ vector of covariables, $\boldsymbol{\alpha} = (\alpha_1, \ldots, \alpha_p)^\top$ corresponds to the vector of regression coefficients, $f(\cdot)$ is unknown smooth arbitrary surface, and $\mathbf{t}_i$ is a two-dimensional covariates vector. To write the model given by Equation (1) in a matrix form, first consider the one-to-one transformation of the vector $\mathbf{f}$ suggested by Green and Silverman [7], stated as

$$\mathbf{f} = \begin{pmatrix} f(\mathbf{t}_1) \\ \vdots \\ f(\mathbf{t}_n) \end{pmatrix} = \mathbf{E}\delta + \mathbf{T}^T \mathbf{a},$$

where $\mathbf{a}$ is a $3 \times 1$ vector with components $a_i$, $\delta$ is a $n \times 1$ vector with components $\delta_i$, $\mathbf{E}$ is a $(n \times n)$ matrix whose elements are given by $E_{ij} = \frac{1}{16\pi} \|\mathbf{t}_i - \mathbf{t}_j\|^2 \log \|\mathbf{t}_i - \mathbf{t}_j\|^2$, with $E_{ii} = 0$ for each $i$, and $\mathbf{T}$ is a $(3 \times n)$ matrix defined as

$$\mathbf{T} = \begin{pmatrix} 1 & 1 & \ldots & 1 \\ \mathbf{t}_1 & \mathbf{t}_2 & \ldots & \mathbf{t}_n \end{pmatrix}.$$

Thus, the Model (1) can be written in a matrix form as

$$\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{E}\delta,$$

where the regression matrix is structured as $X = \begin{pmatrix} \mathbf{T}^T & \mathbf{W} \end{pmatrix}^T$, with $\mathbf{W} = \begin{pmatrix} \mathbf{w}_1^\top, & \ldots, & \mathbf{w}_n^\top \end{pmatrix}^T$, and the vector of regression coefficients as $\boldsymbol{\beta} = \begin{pmatrix} \mathbf{a}^T & \boldsymbol{\alpha} \end{pmatrix}^T = \begin{pmatrix} \beta_1, & \ldots, & \beta_{p+3} \end{pmatrix}$, where

$\beta_j = a_j$ ($j = 1, 2, 3$) and $\beta_j = \alpha_{j-3}$ ($j = 4, \ldots, p+3$); see [9]. Note that this matrix representation of the linear predictor allows us to treat the TPS-GLM as a semiparametric generalized linear model, in which the term $\mathbf{X}\boldsymbol{\beta}$ represents the parametric component and $\mathbf{E}\boldsymbol{\delta}$ the nonparametric component. One of the advantages of the TPS-GLM, apart from being able to model both discrete and continuous variables that belong to the exponential family, is its flexibility to model the non-linear joint effect of covariates through the surface $f$ present in the linear predictor $\eta$. In the context of spatial data, this models allows the effect of coordinates to the incorporated into the modeling process. It is important to note that when the surface $f$ is not present in the linear predictor $\eta$, the model reduces to the classical generalized linear model. However, if the vector $\mathbf{t}$ reduces to a scalar, t, the model reduces to the semiparametric generalized linear model discussed, for instance, by Green and Silverman [7].

*2.2. Penalized Function*

Under the TPS-GLM, we have that $\boldsymbol{\theta} = (\boldsymbol{\beta}^{\top}, \boldsymbol{\delta}^{\top}, \phi)^{\top} \subseteq \mathcal{R}^{p^*}$, with $p^* = (p+3) + n + 1$ parameters. Then, the log-likelihood function is given by

$$L(\boldsymbol{\theta}) = \sum_{i=1}^{n} L_i(\boldsymbol{\theta}), \tag{2}$$

where

$$L_i(\boldsymbol{\theta}) = \left[ \frac{y_i \theta_i - \psi(\theta_i)}{a_i(\phi)} + c(y_i, \phi) \right].$$

To ensure the identifiability of the parameter vector $\boldsymbol{\alpha}$, we assume that $f$ belongs to the function space where all partial derivatives of total order $m$ reside within the Hilbert space $\mathcal{L}^2[E^d]$, the space of square-integrable functions on Euclidean $d$-space. Incorporating a penalty function over $f$, we have that the penalized log-likelihood function can be expressed as (see, for instance, Green and Silverman [7])

$$L_{\mathrm{p}}(\boldsymbol{\theta}, \lambda_f) = L(\boldsymbol{\theta}) + \lambda_f^* J_m^d(f), \tag{3}$$

where $J_m^d(f)$ is a penalty functional measuring the wiggliness of $f$, and $\lambda_f^*(\lambda_f)$ is a constant that depends on the smoothing parameter $\lambda_f \geq 0$. In general, a measure of the curvature of $f$ corresponds to its squared norm, $\|f\|$, defined as

$$J_m^d(f) = \|f\| = \sum_{v_1 + \ldots + v_d = m} \frac{m!}{v_1! \ldots v_d!} \int_{-\infty}^{+\infty} \cdots \int_{-\infty}^{+\infty} \left( \frac{\partial^m f}{\partial t_1^{\alpha_1} \ldots \partial t_d^{\alpha_d}} \right)^2 \prod_{j=1}^{d} \mathrm{dt}_j.$$

For simplicity, in this work, we will consider the case in which $d = 2, m = 2$ and $g = g(t_1, t_2)$. Consequently, the penalty function $J_2^2(f)$ is expressed in the form

$$J_2^2(f) = \int\int_{\mathcal{R}^2} \left\{ \left( \frac{\partial^2 f}{\partial t_1^2} \right)^2 + 2 \left( \frac{\partial^2 f}{\partial t_1 \partial t_2} \right)^2 + \left( \frac{\partial^2 f}{\partial t_2^2} \right)^2 \right\} \mathrm{dt}_1 \mathrm{dt}_2,$$

and measures the rapid variation in $f$ and the departure from local linearity. In this case, the estimation of $f$ leads to a natural thin-plate spline. According to Green and Silverman [7], we may express the penalty functional as $J_2^2(f) = \boldsymbol{\delta}^T \mathbf{E} \boldsymbol{\delta}$. Then, if we consider $\lambda_f^* = -\lambda_f/2$, the penalized log-likelihood function (3) can be expressed as

$$L_{\mathrm{p}}(\boldsymbol{\theta}, \lambda_f) = L(\boldsymbol{\theta}) - \frac{\lambda_f}{2} \boldsymbol{\delta}^T \mathbf{E} \boldsymbol{\delta}. \tag{4}$$

The first term in the right-hand side of Equation (4) measures the goodness-of-fit, while the second terms penalizes the roughness of $f$ with a fixed parameter $\lambda_f$. Selecting appropriate parameters is crucial in the estimation process, as they determine the balance between the goodness-of-fit and the smoothness (or regularity) of the estimated function. It is important to emphasize that selecting appropriate parameters is crucial in the estimation process because they control the trade-off between goodness-of-fit and the smoothness (or regularity) of the estimated function. In this work, the smoothing parameter is selected through the Akaike Criterion (AIC) based on the penalized log-likelihood function given in Equation (3). More details of the method are given in Section 3.7.

## 3. Estimation and Inference

In this section, we discuss the problem of estimating the parameters under the TPS-GLM. Specifically, we derive a weighted iterative process based on the backfitting algorithm and estimate the variance–covariance matrix of our estimator from the penalized Fisher information matrix (see Green [32] and Green and Silverman [7]). A brief discussion of the smoothing parameter selection is also presented.

### 3.1. Penalized Score Function

First, we are going to assume that the function $L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)$ is regular in the sense that it admits first and second partial derivatives with respect to the elements of the parameter vector $\boldsymbol{\theta}$. To obtain the score function for $\boldsymbol{\beta}$, we must calculate $\partial L_{\mathrm{P}_i}(\boldsymbol{\theta}, \lambda_f)/\partial\beta_j$ for $i \in \{1,\ldots,n\}$ and $j \in \{1,\ldots,p+2\}$. After performing some partial derivative operations, we have that the score function for $\boldsymbol{\beta}$ can be written in matrix as follows:

$$\mathbf{U}_{\mathrm{P}}^{\beta}(\boldsymbol{\theta}) = \frac{\partial L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)}{\partial\boldsymbol{\beta}} = \mathbf{X}^{\top}\widetilde{\mathbf{T}}(\mathbf{y} - \boldsymbol{\mu}),$$

where $\mathbf{X}$ is an $(n \times 3 + p)$ matrix whose $i$th row is $\mathbf{x}_i^{\top}$, $\widetilde{\mathbf{T}} = \mathrm{diag}\left[(\mathrm{a}_i(\phi))^{-1}\left(\frac{\partial\mu_i}{\partial\eta_i}\right)V_i^{-1}\right]$ is a $(n \times 3 + p)$ matrix, with $V_i = V(\mu_i) = \partial^2\psi(\theta_i)/\partial\theta_i^2$ the variance function, $\mathrm{a}_i(\phi)$ a function of $\phi$, $\mathbf{y} = (y_1, ..., y_n)^{\top}$ and $\boldsymbol{\mu} = (\mu_1, ..., \mu_n)^{\top}$ are $(n \times 1)$ vectors.

Conversely, to derive the score function for $\delta$, we need to compute $\partial L_{\mathrm{P}_i}(\boldsymbol{\theta}, \lambda_f)/\partial\delta_\ell$ for $i \in \{1,\ldots,n\}$ and $\ell \in \{1,\ldots,n\}$. Again, after some algebraic operations, the score function for $\delta$ can be written in matrix as follows:

$$\mathbf{U}_{\mathrm{P}}^{\delta}(\boldsymbol{\theta}) = \frac{\partial L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)}{\partial\delta} = \mathbf{E}^{\top}\widetilde{\mathbf{T}}(\mathbf{y} - \boldsymbol{\mu}) - \lambda_f\mathbf{E}\delta,$$

where the matrix $\mathbf{E}$ is defined in Section 2.1. Finally, the score function for $\phi$ is given by

$$\mathbf{U}_{\mathrm{P}}^{\phi}(\boldsymbol{\theta}) = \frac{\partial L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)}{\partial\phi} = -\sum_{i=1}^{n}(\mathrm{a}_i(\phi))^{-2}\{y_i\theta_i - \psi(\theta_i)\} + \sum_{i=1}^{n}c'(y_i, \phi),$$

with $c'(y_i, \phi) = \partial c(y_i, \phi)/\partial\phi$, for $i \in \{1,\ldots,n\}$. Thus, the vector of penalized score functions of $\boldsymbol{\theta}$ can be expressed compactly as

$$\mathbf{U}_{\mathrm{P}}(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{U}_{\mathrm{P}}^{\beta}(\boldsymbol{\theta}) \\ \mathbf{U}_{\mathrm{P}}^{\delta}(\boldsymbol{\theta}) \\ \mathbf{U}_{\mathrm{P}}^{\phi}(\boldsymbol{\theta}) \end{pmatrix}.$$

Note that if the model under consideration only considers the parametric component in the linear predictor, that is, the nonparametric component is omitted, the expressions of the remaining score functions are reduced to those obtained under the classical generalized linear model.

### 3.2. Penalized Hessian Matrix

To obtain the penalized Hessian matrix, we must compute the second-derivate of $L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)$ with respect to each element of $\boldsymbol{\theta}$, that is, $\partial^2 L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)/\partial \theta_{j^*} \theta_{\ell^*}$, for $j^*, \ell^* \in \{1, \ldots, p^*\}$. After some algebraic operations, we have that the diagonal elements (block matrices) of the Hessian matrix are given by

$$
\begin{aligned}
\mathbf{L}_{\mathrm{P}}^{\beta\beta} = \frac{\partial^2 L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)}{\partial \boldsymbol{\beta} \partial \boldsymbol{\alpha}^{\top}} &= -\mathbf{X}^{\top} \mathbf{M}^* \mathbf{X}, \\
\mathbf{L}_{\mathrm{P}}^{\delta\delta} = \frac{\partial^2 L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)}{\partial \boldsymbol{\delta} \partial \boldsymbol{\delta}^{\top}} &= -\mathbf{E}^{\top} \mathbf{M}^* \mathbf{E} - \lambda_f \mathbf{E} \qquad \text{and} \\
\mathbf{L}_{\mathrm{P}}^{\phi\phi} = \frac{\partial^2 L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)}{\partial \phi^2} &= \sum_{i=1}^{n} 2(a_i(\phi))^{-3}(y_i\theta_i - \psi(\theta_i)) + \sum_{i=1}^{n} c''(y_i, \phi)),
\end{aligned}
$$

where $\mathbf{M}^* = \mathrm{diag}_{1 \leq i \leq n}\left[(a_i(\phi))^{-1}(\partial \mu_i/\partial \eta_i)^2 V_i^{-1}\right]$ and $c''(y_i, \phi) = \partial^2 c(y_i, \phi)/\partial \phi^2$, for $1 \leq i \leq n$. The elements outside the main diagonal of the Hessian matrix take the form

$$
\begin{aligned}
\mathbf{L}_{\mathrm{P}}^{\beta\delta} = \frac{\partial^2 L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)}{\partial \boldsymbol{\beta} \partial \boldsymbol{\delta}^{\top}} &= -\mathbf{X}^{\top} \mathbf{M}^* \mathbf{E}, \\
\mathbf{L}_{\mathrm{P}}^{\beta_j\phi} = \frac{\partial^2 L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)}{\partial \alpha_j \partial \phi} &= -\sum_{i=1}^{n} (a_i(\phi))^{-2}\left\{(y_i - \mu_i) V_i^{-1} \frac{\partial \mu_i}{\partial \eta_i} x_{ij}\right\} \qquad \text{and} \\
\mathbf{L}_{\mathrm{P}}^{\delta_\ell\phi} = \frac{\partial^2 L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)}{\partial \delta_\ell \partial \phi} &= -\sum_{i=1}^{n} (a_i(\phi))^{-2}\left\{(y_i - \mu_i) V_i^{-1} \frac{\partial \mu_i}{\partial \eta_i} e_{i\ell}\right\},
\end{aligned}
$$

where $x_{ij}$ denotes the $(i, j)$th element of the matrix $\mathbf{X}$ and $e_{ij}$ denotes the $(i, \ell)$th element of the matrix $\mathbf{E}$, for $i \in \{1, \ldots, n\}$, $j \in \{1, \ldots, p+2\}$ and $\ell \in \{1, \ldots, n\}$. Thus, the penalized Hessian matrix can be represented as

$$
\mathbf{L}_{\mathrm{P}}(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{L}_{\mathrm{P}}^{\beta\beta} & \mathbf{L}_{\mathrm{P}}^{\beta\delta} & \mathbf{L}_{\mathrm{P}}^{\beta\phi} \\ \mathbf{L}_{\mathrm{P}}^{\beta\delta^{\top}} & \mathbf{L}_{\mathrm{P}}^{\delta\delta} & \mathbf{L}_{\mathrm{P}}^{\delta\phi} \\ \mathbf{L}_{\mathrm{P}}^{\beta\phi^{\top}} & \mathbf{L}_{\mathrm{P}}^{\delta\phi^{\top}} & \mathbf{L}_{\mathrm{P}}^{\phi\phi} \end{pmatrix}.
$$

It is noteworthy that this matrix simplifies to the Hessian matrix used in generalized linear models when the nonparametric component is absent. The primary application of this matrix lies in the normal curvature, which is essential for developing the local influence technique. This will be discussed in the following section.

### 3.3. Penalized Expected Information Matrix

By taking the expectation of the matrix $-\mathbf{L}_{\mathrm{P}}(\boldsymbol{\theta})$, we derive the penalized expected information matrix, which is of dimension $(p^* \times p^*)$, as follows:

$$
\mathcal{J}_{\mathrm{P}}(\boldsymbol{\theta}) = -\mathrm{E}\left[\frac{\partial^2 L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T}\right].
$$

This matrix assumes the following diagonal structure in blocks:

$$
\mathcal{J}_{\mathrm{P}}(\boldsymbol{\theta}) = \begin{pmatrix} \mathcal{J}_{\mathrm{P}}^{\beta\delta}(\boldsymbol{\theta}) & \mathbf{0} \\ \mathbf{0} & \mathcal{J}_{\mathrm{P}}^{\phi\phi}(\boldsymbol{\theta}) \end{pmatrix},
$$

where

$$\mathfrak{J}_{\mathrm{P}}^{\beta\delta}(\boldsymbol{\theta}) \quad = \quad \begin{pmatrix} \mathbf{X}^\top \mathbf{M}^* \mathbf{X} & \mathbf{X}^\top \mathbf{M}^* \mathbf{E} \\ \mathbf{E}^\top \mathbf{M}^* \mathbf{X} & \mathbf{E}^\top \mathbf{M}^* \mathbf{E} + \lambda_f \mathbf{E} \end{pmatrix}$$

and

$$\mathfrak{J}_{\mathrm{P}}^{\phi\phi}(\boldsymbol{\theta}) \quad = \quad \sum_{i=1}^{n} -2(\mathrm{a}_i(\phi))^{-3}(\mu_i \theta_i - \psi(\theta_i)) - \sum_{i=1}^{n} \mathrm{E}(c''(y_i, \phi)),$$

with $c''(y_i, \phi) = \partial^2 c(y_i, \phi)/\partial \phi^2$ for $i \in \{1, \dots, n\}$.

*3.4. Derivation of the Iterative Process*

The value of $\boldsymbol{\theta}$ that maximizes $L_{\mathrm{P}}(\boldsymbol{\theta}, \lambda_f)$, called maximum penalized likelihood estimate (MPLE) and denoted by $\widehat{\boldsymbol{\theta}}$, is carried out by solving the corresponding estimation equations. Let $\boldsymbol{\theta} = \begin{pmatrix} \boldsymbol{\theta}_1^\top & \theta_2 \end{pmatrix}^\top$, where $\boldsymbol{\theta}_1 = \begin{pmatrix} \boldsymbol{\beta}^\top & \boldsymbol{\delta}^\top \end{pmatrix}^\top$ and $\theta_2 = \phi$. In addition, consider the partition of the score function vector $\mathbf{U}_{\mathrm{P}}(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{U}_{\mathrm{P}}^{1^\top}(\boldsymbol{\theta}) & \mathbf{U}_{\mathrm{P}}^{2^\top}(\boldsymbol{\theta}) \end{pmatrix}^\top$, where $\mathbf{U}_{\mathrm{P}}^{1^\top}(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{U}_{\mathrm{P}}^{\beta^\top}(\boldsymbol{\theta}) & \mathbf{U}_{\mathrm{P}}^{\delta^\top}(\boldsymbol{\theta}) \end{pmatrix}$ and $\mathbf{U}_{\mathrm{P}}^2(\boldsymbol{\theta}) = \mathbf{U}_{\mathrm{P}}^\phi(\boldsymbol{\theta})$. In order to estimate $\boldsymbol{\theta}$ based on penalized likelihood function given by Equation (4), we have to solve the equations

$$\begin{cases} \mathbf{U}_{\mathrm{P}}^1(\boldsymbol{\theta}) & = & \mathbf{0} \\ \mathbf{U}_{\mathrm{P}}^2(\boldsymbol{\theta}) & = & \mathbf{0} \,. \end{cases}$$

These estimating equations are nonlinear, and necessitate an iterative approach for their solution. An alternative frequently proposed in the context of generalized linear models is the Fisher scoring algorithm (Nelder and Wedderburn, [33]), considering the fact that in some situations the matrix $-\mathbf{L}_{\mathrm{P}}(\boldsymbol{\theta})$ can be non-positive definite. Then, the algorithm for estimating $\boldsymbol{\theta}_1$, with $\phi$ fixed, is given by

$$\boldsymbol{\theta}_1^{\mathrm{new}} = \boldsymbol{\theta}_1^{\mathrm{old}} + (\mathfrak{J}_{\mathrm{P}}^{\beta\delta}(\boldsymbol{\theta})^{-1})^{\mathrm{old}} \mathbf{U}_{\mathrm{P}}^{1^{\mathrm{old}}}(\boldsymbol{\theta}),$$

which is equivalent to solving the matrix equation

$$\begin{pmatrix} \mathbf{I} & \mathbf{S}_\beta^{\mathrm{old}} \mathbf{E} \\ \mathbf{S}_\delta^{\mathrm{old}} \mathbf{X} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta}^{\mathrm{new}} \\ \boldsymbol{\delta}^{\mathrm{new}} \end{pmatrix} = \begin{pmatrix} \mathbf{S}_0^{\mathrm{old}} \mathbf{z}^{\mathrm{old}} \\ \mathbf{S}_1^{\mathrm{old}} \mathbf{z}^{\mathrm{old}} \end{pmatrix}, \tag{5}$$

where $\mathbf{z}^{\mathrm{old}} = (\mathbf{y} - \boldsymbol{\mu}^{\mathrm{old}}) + \boldsymbol{\eta}^{\mathrm{old}}$, with $\mathbf{S}_\vartheta^{\mathrm{old}}$ defined as

$$\mathbf{S}_\vartheta \quad = \quad \begin{cases} (\mathbf{X}^\top \mathbf{M}^{*\mathrm{old}} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{M}^{*\mathrm{old}} & \vartheta = \beta \\[2ex] (\mathbf{E}^\top \mathbf{M}^{*\mathrm{old}} \mathbf{E} + \lambda_f \mathbf{E})^{-1} \mathbf{E}^\top \mathbf{M}^{*\mathrm{old}} & \vartheta = \delta \,. \end{cases}$$

Consequently, the weighted back-fitting (Gauss–Seidel) iterations for simultaneously updating $\boldsymbol{\beta}$ and $\boldsymbol{\delta}$ are given by

$$\boldsymbol{\beta}^{\mathrm{new}} \quad = \quad \mathbf{S}_\beta^{\mathrm{old}} (\mathbf{z}^{\mathrm{old}} - \mathbf{E}\boldsymbol{\delta}^{\mathrm{old}}), \tag{6}$$

$$\boldsymbol{\delta}^{\mathrm{new}} \quad = \quad \mathbf{S}_\delta^{\mathrm{old}} (\mathbf{z}^{\mathrm{old}} - \mathbf{X}\boldsymbol{\beta}^{\mathrm{old}}), \tag{7}$$

It is crucial to note that the system of Equations (5) is consistent, and the back-fitting algorithm converges to a solution for any initial values, provided that the weight matrix

$\mathbf{M}^*$ is symmetric and positive definite. Additionally, if the parametric component $\mathbf{w}_i^\top \boldsymbol{\beta}$ is absent in the linear predictor, the estimator of $\delta$ is given by:

$$\delta^{\text{new}} = \mathbf{S}_\delta^{\text{old}} \mathbf{z}^{\text{old}} .$$

The MPLE of the dispersion parameter, $\widehat{\theta}_2 = \widehat{\phi}$, can be determined through the following iterative procedure:

$$\theta_2^{\text{new}} = \theta_2^{\text{old}} - (\mathcal{J}_{\text{P}}^{\phi\phi}(\boldsymbol{\theta})^{-1})^{\text{old}} \mathbf{U}_{\text{P}}^{2^{\text{old}}}(\boldsymbol{\theta}) .$$

Summarizing, each iteration of the Fisher scoring algorithm updates $\boldsymbol{\beta}$ and $\delta$ using Equations (6) and (7), and evaluating matrices $\mathbf{S}_\vartheta$ and $\mathbf{M}^*$ at the MPLE of $\boldsymbol{\theta}$ obtained in the previous iteration, that is, $\boldsymbol{\theta}^{\text{old}}$, until convergence is obtained. The joint iterative process that resolves $\mathbf{U}_{\text{P}}(\boldsymbol{\theta}) = \mathbf{0}$ is presented below.

### 3.5. Estimation of Surface

To obtain the MPLE of $\mathbf{f}$, we must consider its one-to-one representation given in Equation (2) and MPLE obtained from the iterative process described above. Indeed, we have that $\widehat{\mathbf{f}}$ can be obtained as

$$\widehat{\mathbf{f}} = \mathbf{E}\widehat{\delta} + \mathbf{T}^T\widehat{\mathbf{a}}, ,$$

where $\widehat{\delta}$ and $\widehat{\mathbf{a}}$ are the MPLE of $\delta$ and $\widehat{\mathbf{a}}$, respectively. Note that vector $\widehat{\mathbf{a}}$ corresponds to the first three elements of vector $\widehat{\boldsymbol{\beta}}$. Consequently, $\widehat{\mathbf{f}}$ is a natural thin-plate spline. Details of the conditions that guarantee this result are given, for example, in Green and Silverman [7] and Wood [34].

### 3.6. Approximate Standard Errors

In this study, we propose approximating the variance–covariance matrix of $\widehat{\boldsymbol{\theta}}$ by using the inverse of the penalized Fisher information matrix. Specifically, we have that

$$\widehat{\text{Cov}}(\widehat{\boldsymbol{\theta}}) \approx \mathcal{J}_{\text{P}}^{-1}(\boldsymbol{\theta})\big|_{\widehat{\boldsymbol{\theta}}} .$$

If we are interested in drawing inferences for $\boldsymbol{\beta}$, the approximate variance–covariance matrix can be estimated by using the corresponding block-diagonal matrix obtained from $\mathcal{J}_{\text{P}}^{-1}(\boldsymbol{\theta})$, similarly for $\mathbf{f}$ and $\phi$.

### 3.7. On Degrees of Freedom and Smoothing Parameter

For the TPS-GLM, the degree of freedom ($df$) associated with the smooth surface is given by (see, for instance Green and Silverman [7])

$$df(\lambda_f) = \text{tr}(\mathbf{E}^\top \mathbf{S}_\delta),$$

which approximately represents the number of effective parameters used in the modeling process to estimate the smooth surface $f$.

Regarding the selection of the smoothing parameter, we propose to use the Akaike Information Criterion (AIC) (see, for instance, [24,35]), defined as

$$AIC(\lambda_f) = -2L_p(\boldsymbol{\theta}, \lambda_f)\big|_{\widehat{\boldsymbol{\theta}}} + 2[1 + p + df(\lambda_f)] ,$$

where $L_p(\boldsymbol{\theta}, \lambda_f)$ denote the penalized likelihood function evaluated at MPLE of $\boldsymbol{\theta}$, and $p$ denote the number of parameters in $\boldsymbol{\beta}$. As usual, the idea is to select the value of $\lambda_f$ that minimizes $AIC(\lambda_f)$.

## 4. Local Influence

In this section, we extend the local influence technique to evaluate the sensitivity of the MPLE under the TPS-GLM. Specifically, we present some theoretical aspects of the method and, subsequently, we derive the normal curvature for three perturbation schemes.

### 4.1. Local Influence Analysis

Consider $\boldsymbol{\omega} = (\omega_1, \ldots, \omega_n)^\top$, an $n \times 1$ vector of perturbations restricted to some open subset $\Omega \subset \mathbb{R}^n$. Let $L_p(\boldsymbol{\theta}, \lambda_f \mid \boldsymbol{\omega})$ denote the logarithm of the perturbed penalized likelihood function. Assume there exists a vector of non-perturbation $\boldsymbol{\omega}_0 \in \Omega$, such that $L_p(\boldsymbol{\theta}, \lambda_f \mid \boldsymbol{\omega}_0) = L_p(\boldsymbol{\theta}, \lambda_f)$. To evaluate the influence of small perturbations on the MPL estimate $\widehat{\boldsymbol{\theta}}$, we can consider the penalized likelihood displacement given by:

$$2[L_p(\widehat{\boldsymbol{\theta}}, \lambda_f) - L_p(\widehat{\boldsymbol{\theta}}_\omega, \lambda_f)] \geq 0,$$

where $\widehat{\boldsymbol{\theta}}_\omega$ is the MPL estimate under $L_p(\boldsymbol{\theta}, \lambda_f \mid \boldsymbol{\omega})$. The measure $\mathrm{LD}(\boldsymbol{\omega})$ is useful for assessing the distance between $\widehat{\boldsymbol{\theta}}$ and $\widehat{\boldsymbol{\theta}}_\omega$. Cook [10] suggested examining the local behavior of $\mathrm{LD}(\boldsymbol{\omega})$ around $\boldsymbol{\omega}_0$. The procedure involves selecting a unit direction $\boldsymbol{d} \in \Omega$, with $\|\boldsymbol{d}\| = 1$, and then plotting $\mathrm{LD}(\boldsymbol{\omega}_0 + a\boldsymbol{d})$ against $a$, where $a \in \mathbb{R}$. This plot, called the lifted line, can be characterized by considering the normal curvature $C_d(\boldsymbol{\theta})$ around $a = 0$. The suggestion is to assume the direction $\boldsymbol{d} = \boldsymbol{d}_{\max}$ corresponding to the largest curvature $C_{d_{\max}}(\boldsymbol{\theta})$. The index plot of $\boldsymbol{d}_{\max}$ can identify those cases that, under small perturbations, have a significant potential influence on $\mathrm{LD}(\boldsymbol{\omega})$. According to Cook [10], the normal curvature in the unit direction $\boldsymbol{d}$ is expressed as

$$C_d(\boldsymbol{\theta}) = -2\{\boldsymbol{d}^\top \boldsymbol{\Delta}_{\mathrm{p}}^\top \mathbf{L}_{\mathrm{p}}^{-1} \boldsymbol{\Delta}_{\mathrm{p}} \boldsymbol{d}\},$$

with

$$\mathbf{L}_{\mathrm{p}} = \left. \frac{\partial^2 L_p(\boldsymbol{\theta}, \lambda_f)}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right|_{\boldsymbol{\theta} = \widehat{\boldsymbol{\theta}}} \quad \text{and} \quad \boldsymbol{\Delta}_{\mathrm{p}} = \left. \frac{\partial^2 L_p(\boldsymbol{\theta}, \lambda_f | \boldsymbol{\omega})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\omega}^\top} \right|_{\boldsymbol{\theta} = \widehat{\boldsymbol{\theta}},\ \boldsymbol{\omega} = \boldsymbol{\omega}_0}.$$

Note that $-\mathbf{L}_{\mathrm{p}}$ represents the penalized observed information matrix evaluated at $\widehat{\boldsymbol{\theta}}$ (see Section 3.2), and $\boldsymbol{\Delta}_{\mathrm{p}}$ is the penalized perturbation matrix evaluated at $\widehat{\boldsymbol{\theta}}$ and $\boldsymbol{\omega}_0$. It is essential to highlight that $C_d(\boldsymbol{\theta})$ denotes the local influence on the estimate $\widehat{\boldsymbol{\theta}}$ after perturbing the model or data. Escobar and Meeker [36] suggested examining the normal curvature in the direction $\boldsymbol{d} = \boldsymbol{e}_i$, where $\boldsymbol{e}_i$ is an $n \times 1$ vector with a one at the $i$th position and zeros elsewhere. Consequently, the normal curvature, referred to as the total local influence of the $i$th case, takes the form $C_{e_i}(\boldsymbol{\theta}) = 2|c_{ii}|$ for $i \in \{1, \ldots, n\}$, where $c_{ii}$ is the $i$th principal diagonal element of the matrix $\boldsymbol{C} = \boldsymbol{\Delta}_{\mathrm{p}}^\top \mathbf{L}_{\mathrm{p}}^{-1} \boldsymbol{\Delta}_{\mathrm{p}}$.

### 4.2. Derivation of the Normal Curvature

Typically, the perturbation schemes used in the analysis of local influence are determined by the structure of the model under consideration, as discussed by Billor and Loynes [37]. These schemes can generally be divided into two main categories: perturbations to the model (to examine changes in assumptions) or perturbations to the data. For instance, we might consider perturbing the response variable or the explanatory variables. The motivation for employing these perturbation schemes often includes addressing issues such as the presence of outliers or the occurrence of measurement errors in the data. Subsequently, we will present the formulas for the matrix $\boldsymbol{\Delta}_{\mathrm{p}}$ for various perturbation schemes.

Consider the weights assigned to the observations in the penalized log-likelihood function, given by:

$$L_p(\boldsymbol{\theta}, \lambda_f | \boldsymbol{\omega}) \;\;=\;\; L(\boldsymbol{\theta} | \boldsymbol{\omega}) - \sum_{i=1}^n \frac{\lambda_f}{2} \boldsymbol{\delta}^\top \mathbf{E} \boldsymbol{\delta} \,,$$

where $L(\boldsymbol{\theta}|\omega) = \sum_{i=1}^{n} \omega_i L_i(\boldsymbol{\theta})$, $\boldsymbol{\omega} = (\omega_1, \ldots, \omega_n)^\top$ is the vector of weights, with $0 \le \omega_i \le 1$. In this case, the vector of no perturbation is given by $\boldsymbol{\omega}_0 = 1_{(n \times 1)}$. Differentiating $L_p(\boldsymbol{\theta}, \lambda_f | \boldsymbol{\omega})$ with respect to the elements of $\boldsymbol{\theta}$ and $\boldsymbol{\omega}^T$, we have that the matrix $\boldsymbol{\Delta}_{\mathrm{p}}$ takes the form

$$\boldsymbol{\Delta}_{\mathrm{p}} = \begin{pmatrix} \mathbf{X}^\top \mathbf{D}_\tau \\ \mathbf{E}^\top \mathbf{D}_\tau \\ \hat{\mathbf{u}}^\top \end{pmatrix},$$

where the matrix $\mathbf{D}_\tau = \mathrm{diag}_{1 \le i \le n}(\tau_i)$ and $\mathbf{u} = (u_1, \ldots, u_n)^\top$, with $\tau_i = (a_i(\phi))^{-1}(y_i - \partial \psi(h(\eta_i))/\partial h(\eta_i))\partial h(\eta_i)/\partial \eta_i$, $h(\eta_i) = \psi'^{-1}(\eta_i)$, $\psi'^{-1}(\cdot)$ denotes the inverse function of $\psi'(\cdot)$, $u_i = -(a_i(\phi))^{-2}(y_i h(\eta_i) - \psi(h(\eta_i)) + c'(y_i, \phi)\mathbf{e}_{in}^\top$, and $\mathbf{e}_{in}$ a vector with 1 at the $i$th position and zero elsewhere.

To perturb the response variable values, we consider $y_{i\omega} = y_i + \omega_i$ for $i \in \{1, \ldots, n\}$, where $\boldsymbol{\omega} = (\omega_1, \ldots, \omega_n)^\top$ is the vector of perturbations. The vector of no perturbation is $\boldsymbol{\omega} = (0, \ldots, 0)^\top$. The perturbed penalized log-likelihood function is constructed from Equation (3) with $y_i$ replaced by $y_{i\omega}$, as follows:

$$L_p(\boldsymbol{\theta}, \lambda_g | \boldsymbol{\omega}) \quad = \quad L(\boldsymbol{\theta}|\boldsymbol{\omega}) - \sum_{i=1}^{n} \frac{\lambda_f}{2} \delta^\top \mathbf{E} \delta \,,$$

where $L(\cdot)$ is defined in Equation (2) with $y_{i\omega}$ replacing $y_i$. By differentiating $L_p(\boldsymbol{\theta}, \lambda \mid \boldsymbol{\omega})$ with respect to the elements of $\boldsymbol{\theta}$ and $\omega_i$, and after some algebraic manipulation, we obtain:

$$\boldsymbol{\Delta}_{\mathrm{p}} = \begin{pmatrix} \mathbf{X}^\top \mathbf{D}_c \\ \mathbf{E}^\top \mathbf{D}_c \\ \hat{\mathbf{d}}^\top \end{pmatrix},$$

where the matrix $\mathbf{D}_c = \mathrm{diag}_{1 \le i \le n}(c_i)$ and $\mathbf{d} = (d_1, \ldots, d_n)^\top$, with $c_i = \partial h(\eta_i)/\partial \eta_i$ and $d_i = -(a_i(\phi))^{-2}(h(\eta_i)\mathbf{e}_{in}^\top + c'(y_{i\omega}, \phi)/\partial \omega_i)$, with $\mathbf{e}_{in}$ denoting a vector with 1 at the $i$th position and zero elsewhere.

## 5. Applications

In this section, we show the applicability of the TPS-GLM and the local influence method with two real data applications. The model estimation and diagnostics have been implemented using `MatLab` 9.13.0 (R2022b) software [38] (the developed code is available on request by the authors).
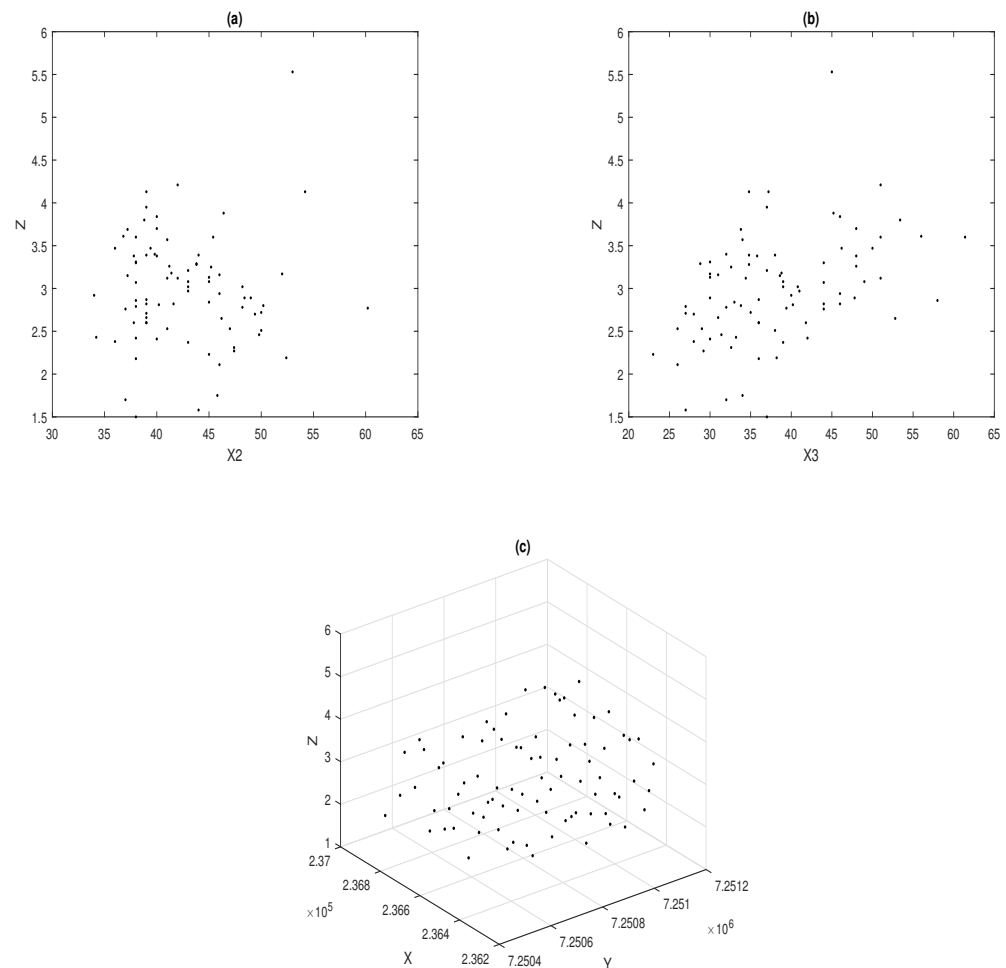
### 5.1. Wypych Data

The first dataset we use to illustrate the applicability of the TPS-GLM consists of 83 sample points within a 46.6-hectare agricultural area in Wypych, located at latitude 24°50′24″ S and longitude 53°36′36″ W, with an average altitude ranging from 589 to 660 m. The data were collected during the 2006/2007 agricultural year in the western region of Paraná State, Brazil (see [39], Appendix 4). The soil is classified as Dystroferric Red Latosol with a clayey texture. The region's climate is mesothermal, super-humid temperate, classified as Cfa according to (Köeppen), with a mean annual temperature of 21 °C. The 83 georeferenced points were determined by a regular grid of $75 \times 75$ m using a global positioning system (GPS). The collected variables were as follows:

- Soya: average of soybean yield (t/ha).
- Height: average height (cm)of plants at the end of the production process.
- Pods: average number of pods.
- Lat: latitude (UTM).
- Long: longitude (UTM).

The original objective was to investigate the spatial variability of soybean yield (Soya) in the studied area based on the covariates: average plant height, average number of pods

per plant, latitude, and longitude. Figure 1 shows the scatterplots between the response variable Soya and the explanatory variables Height and Pods. In addition, the plot of the response variable against the coordinates is shown. Clearly, from Figure 1a,b, it can be seen that the explanatory variables Height (X2) and Pods (X3) are linearly related to the response variable Soya (Z). The spatial effect given by the coordinates (X,Y) will be incorporated into the model through a smooth surface.



**Figure 1.** Scatter plots: Soya versus Height (**a**), Soya versus Pods (**b**), and Soya versus coordinates in UTM (**c**).

### 5.1.1. Fitting the TPS-GLM

Based on the above analysis, we propose the TPS-GLM, introduced in Section 2, to model the trends present in the Wypych data. Specifically, we are going to assume that the response variable Soya belongs to the Gaussian family, and that the link function is the identity. Therefore, the model is expressed as follows:

$$g(\mu_i) = \mu_i = \beta_0 + \beta_1 \text{Height}_i + \beta_2 \text{Pods}_i + f(\text{Lat}_i, \text{Long}_i) \qquad i \in \{1, \ldots, 83\},$$

where $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2)^\top$ correspond to the regression coefficients associated with the parametric component of the model, and $f(\cdot)$ is a smooth surface. Table 1 lists the MPLE of $\boldsymbol{\beta}$. The respective asymptotic standard errors are presented in parentheses.

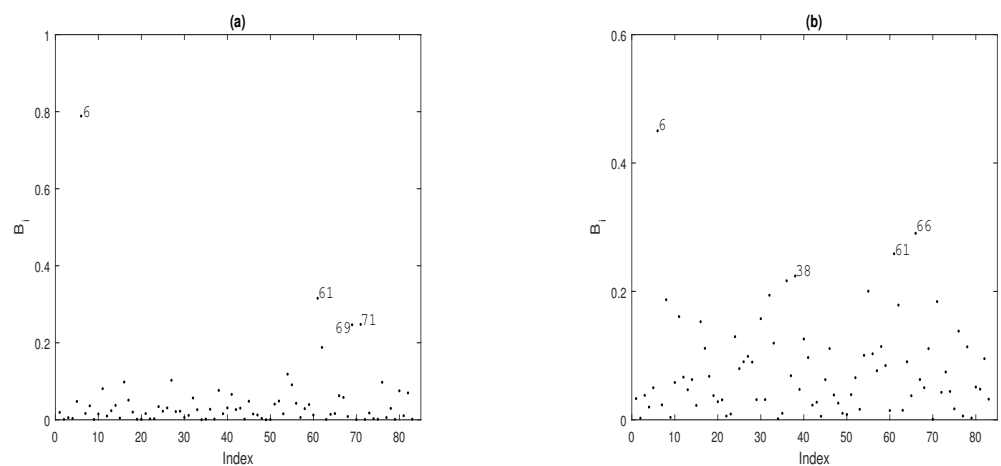**Table 1.** MPLEs with their standard errors (within parenthesis), AIC and R$^2$(Adj).

| | Model | |
|---|---|---|
| **Parameters** | **Gaussian Linear** | **TPS-GLM** |
| $\beta_0$ | 1.1921 (0.672) | 0.497 (0.751) |
| $\beta_1$ | 0.0116 (0.0128) | 0.032 (0.015) |
| $\beta_2$ | 0.0339 (0.0079) | 0.030 (0.008) |
| AIC | 149.99 | 139.9992 |
| R$^2$(Adj) | 0.168 | 0.315 |

The value of the smoothing parameter $\lambda_f$ was selected in such a way that the AIC criterion was minimized. The adjusted determinant coefficients (R$^2$(Adj)) are evaluated for assessing the goodness-of-fit of the two models. It is important to note that our model have a lower AIC and an higher R$^2$(Adj), compared to the multiple regression model that does not consider the spatial effect. Figure 2a shows the QQ-plot for the standardized residuals, whose adjustment to the Gaussian TPS-GLM seems to be reasonable. However, the presence of some atypical observations is observed in one of the tails of the distribution. Figure 2b displays the scatter plot between the observed values, Soya, and their estimated values, $\widehat{\text{Soya}}$. Considering the trend of the points, we conclude that the estimates are good, since they generate consistent adjusted values of the response variable.



**Figure 2.** QQ-plots of the standardized residuals for the TPS-GLM with its confidence interval (dashed lines) (**a**) and scatterplot between Soya and $\widehat{\text{Soya}}$ (**b**), under model fitted to Wypych data.

### 5.1.2. Diagnostic Analysis

To identify potentially influential observations on the MPLE under the fitted Gaussian TPS-GLM for the Wypych data, we present several index plots of $B_i = B_{e_i}(\gamma)$ for $\gamma = \beta, \delta$. Figure 3 shows the index plot $B_i$ for the case-weight perturbation scheme under the fitted model. Figure 3a reveal that the observations #6, #61, #69 and #71 are more influential on $\widehat{\beta}$, whereas the observations #6, #66, #61 and #38 are more influential on $\widehat{\delta}$; see Figure 3b. When we perturb the response variable additively, we have that the observations #80, #32, #75 and #88 are more influential on $\widehat{\beta}$; see Figure 4a. Regarding $\widehat{\delta}$, observations #3, #42 and #80 appear as slightly influential as seem in Figure 4b.

**Figure 3.** Index plots of $B_i$ for assessing local influence on $\widehat{\beta}$ (**a**) and $\widehat{\delta}$ (**b**), considering case-weight perturbation.



**Figure 4.** Index plots of $B_i$ for assessing local influence on $\widehat{\beta}$ (**a**) and $\widehat{\delta}$ (**b**), considering response variable additive perturbation.

We conclude that the maximum penalized likelihood estimates (MPLE) of the regression coefficients and the smooth surface exhibit sensitivity to modifications made to the data or the model. This analysis has shown that observations identified as influential for the parametric component do not necessarily exert influence on the non-parametric component, and vice versa. For instance, under the case-weight perturbation scheme, observations #69 and #71 were detected as influential for the parametric component, but not for the nonparametric component.

### 5.1.3. Confirmatory Analysis

Table 2 displays the relative changes experienced by each element in the vector of regression coefficients. In this analysis, we only consider the three most influential observations under the case-weight perturbation scheme. As can be seen in this table, observations #6, #61 and #69 generate significant changes in the estimates. Still, no relevant inferential changes were noted. However, the AIC and &$R^2$(Adj) present some differences once the above observations are dropped.

**Table 2.** Relative changes (RCs) (in %) in the MPL estimates of $\beta_j$ in cases-weight perturbation under the TPS-GLM. The last two columns indicate the AIC and $R^2$(Adj) of the model with dropped observations.

| Dropped Obs. | Parameters and Relatives Changes | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $\beta_0$ | $\beta_1$ | $\beta_2$ | $RC_{\beta_0}$ | $RC_{\beta_1}$ | $RC_{\beta_2}$ | AIC | $R^2$(Adj) |
| 6 | 1.696 | 0.009 | 0.023 | 122.59 | 63.93 | 27.52 | 125.50 | 0.267 |
| 61 | 0.987 | 0.022 | 0.027 | 29.47 | 15.19 | 12.23 | 131.96 | 0.356 |
| 69 | 0.432 | 0.035 | 0.028 | 43.33 | 35.04 | 10.85 | 138.07 | 0.326 |
| 6-61 | 1.996 | 0.002 | 0.022 | 161.89 | 92.06 | 27.65 | 116.20 | 0.218 |
| 6-69 | 1.481 | 0.014 | 0.023 | 94.35 | 45.61 | 26.09 | 124.52 | 0.268 |
| 61-69 | 0.704 | 0.029 | 0.028 | 7.59 | 9.58 | 10.93 | 131.03 | 0.358 |
| 6-61-69 | 1.868 | 0.005 | 0.023 | 145.16 | 81.03 | 26.90 | 115.83 | 0.310 |

On the other hand, Table 3 shows the relative changes in the vector of regression coefficients under the additive perturbation scheme of the response variable. Here, we consider the four most influential observations. As can be seen in the table, observations #32, #69, #75 and #80 generate important relative changes in the estimates of the parametric component of the model. However, no significant inferential changes were observed. About the AIC and &$R^2$(Adj), there are not evident differences.

**Table 3.** Relative changes (RCs) (in %) in the MPL estimates of $\beta_j$ in response variable perturbation under the TPS-GLM. The last two columns indicate the AIC and $R^2$(Adj) of the model with dropped observations.
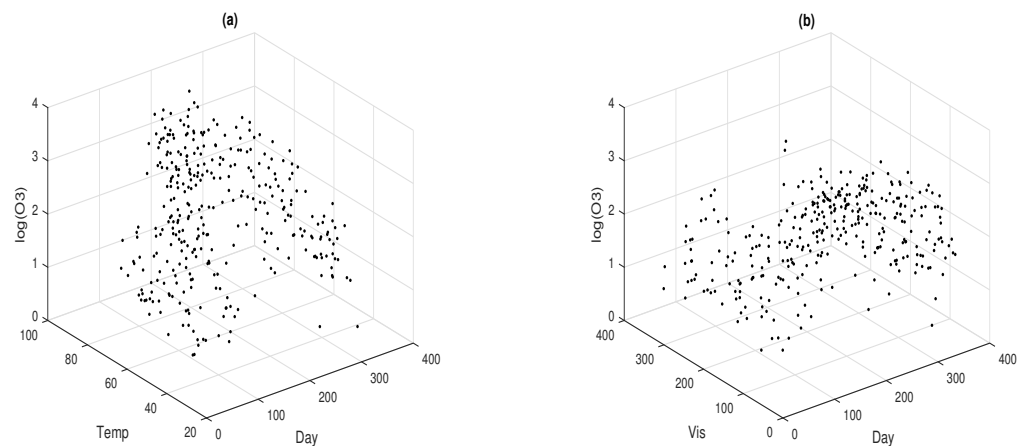
| Dropped Obs. | Parameters and Relatives Changes | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $\beta_0$ | $\beta_1$ | $\beta_2$ | $RC_{\beta_0}$ | $RC_{\beta_1}$ | $RC_{\beta_2}$ | AIC | $R^2$(Adj) |
| 32 | 0.701 | 0.027 | 0.029 | 8.000 | 4.62 | 5.23 | 138.73 | 0.319 |
| 69 | 0.432 | 0.034 | 0.027 | 43.33 | 29.0 | 12.22 | 138.07 | 0.326 |
| 75 | 0.760 | 0.028 | 0.027 | 0.24 | 6.22 | 12.65 | 139.44 | 0.307 |
| 80 | 0.699 | 0.028 | 0.029 | 8.21 | 7.86 | 8.16 | 139.01 | 0.311 |
| 32-69 | 0.382 | 0.035 | 0.030 | 49.87 | 32.82 | 4.00 | 136.81 | 0.33 |
| 32-75 | 0.700 | 0.027 | 0.030 | 8.17 | 4.12 | 4.90 | 138.11 | 0.319 |
| 32-80 | 0.621 | 0.028 | 0.031 | 18.49 | 6.34 | 0.16 | 137.63 | 0.316 |
| 69-75 | 0.430 | 0.035 | 0.028 | 43.61 | 34.96 | 10.96 | 137.53 | 0.318 |
| 69-80 | 0.333 | 0.036 | 0.029 | 56.33 | 38.32 | 5.39 | 136.99 | 0.322 |
| 75-80 | 0.695 | 0.028 | 0.029 | 8.85 | 7.25 | 7.90 | 138.39 | 0.302 |
| 32-69-75 | 0.381 | 0.035 | 0.030 | 49.98 | 32.33 | 3.74 | 136.22 | 0.322 |
| 32-75-80 | 0.621 | 0.028 | 0.031 | 18.54 | 5.23 | 0.96 | 136.93 | 0.308 |
| 69-75-80 | 0.333 | 0.036 | 0.029 | 56.26 | 37.40 | 5.09 | 136.42 | 0.314 |
| 32-69-75-80 | 0.271 | 0.035 | 0.032 | 64.47 | 30.15 | 3.22 | 134.96 | 0.320 |

*5.2. Ozone Concentration Data*

For our analysis, we utilize data from a study examining the relationship between atmospheric ozone concentration (O3) and various meteorological variables in the Los Angeles Basin for a sample of 330 days in 1976. The data were initially presented by Breiman and Friedman [40], and are available for download from various public repositories. Although the dataset includes several variables, in this application, we will consider only three explanatory variables, which are detailed in the following.

- O3: daily maximum one-hour average ozone concentration in Upland, CA, measured in parts per million (ppm).
- Temp: Sandburg Air Base temperature, in Celsius.
- Vis: visibility, in miles.
- Day: calendar day.

Figure 5 contains the dispersion graphs between the outcome variable (log(O3)) and each one of the explanatory variables Temp, Vis and Day.

**Figure 5.** 3D plots between the response variable and the explanatory variables: logarithm of ozone data versus temperature and day variables (**a**), and logarithm of ozone data versus visibility and day variables (**b**).

Figure 5a shows a curved surface in the relationship between the variable log(O3) and the joint effect of the explanatory variables Temp and Day, whereas the relationship between log(O3) and the joint effect of the explanatory variable Vis and Day shows less curve; see Figure 5b. This graphical analysis recommends the inclusion in the model of a nonparametric component, specifically a surface, that can explain the relationship between log(O3) and the combined effect of the explanatory variables Temp and Day. For simplicity, in this work, we will include the effect of the explanatory variable Vis in a linear form. To begin our analysis, we are going to consider the fit of a GLM assuming that the variable of interest O3 is Poisson distributed with mean $\mu_i$ and logarithmic link function. Different structures of the linear predictor for the explanatory variables Vis, Temp, and Day will be considered (see Table 4).

**Table 4.** Four structures of the linear predictor for the explanatory variables Vis, Temp, and Day, assuming that the response variable log(O3) follows a POISSON($\mu_i$) distribution.

| Model | $g(\mu_i) = \log(\mu_i)$ |
|:-----:|:--------------------------|
| I | $\beta_0 + \beta_1 \text{Vis}_i + \beta_2 \text{Temp}_i + \beta_3 \text{Day}_i$ |
| II | $\beta_0 + \beta_1 \text{Vis}_i + \beta_2 \text{Temp}_i + f(\text{Day}_i)$ |
| III | $\beta_0 + \beta_1 \text{Vis}_i + \beta_2 \text{Temp}_i + \beta_3 \text{Day}_i + \beta_4 \text{Temp}_i \times \text{Day}_i$ |
| IV | $\beta_0 + \beta_1 \text{Vis}_i + f(\text{Temp}_i, \text{Day}_i)$ |

For Model I, we consider only the individual effects of the explanatory variables Vis, Temp and Day. Note that all these effects were incorporated in a linear form in the systematic component of the model. For Model II, we consider the inclusion of a nonparametric term to model the nonlinear effects of the explanatory variable Day; see Ibacache et al. [41]. Model III considers a systematic component that contains the individual effects of the explanatory variables Vis, Temp and Day, in addition to the incorporation of the interaction effect between the explanatory variables Temp and Day. Here, the interaction effect is introduced linearly in the model. Model IV corresponds to a TPS-GLM where the joint effect of the Temp and Day explanatory variables is included nonlinearly by using smooth surface. Table 5 contains the ML and MPL estimates associated with the parametric component for the four fitted models.
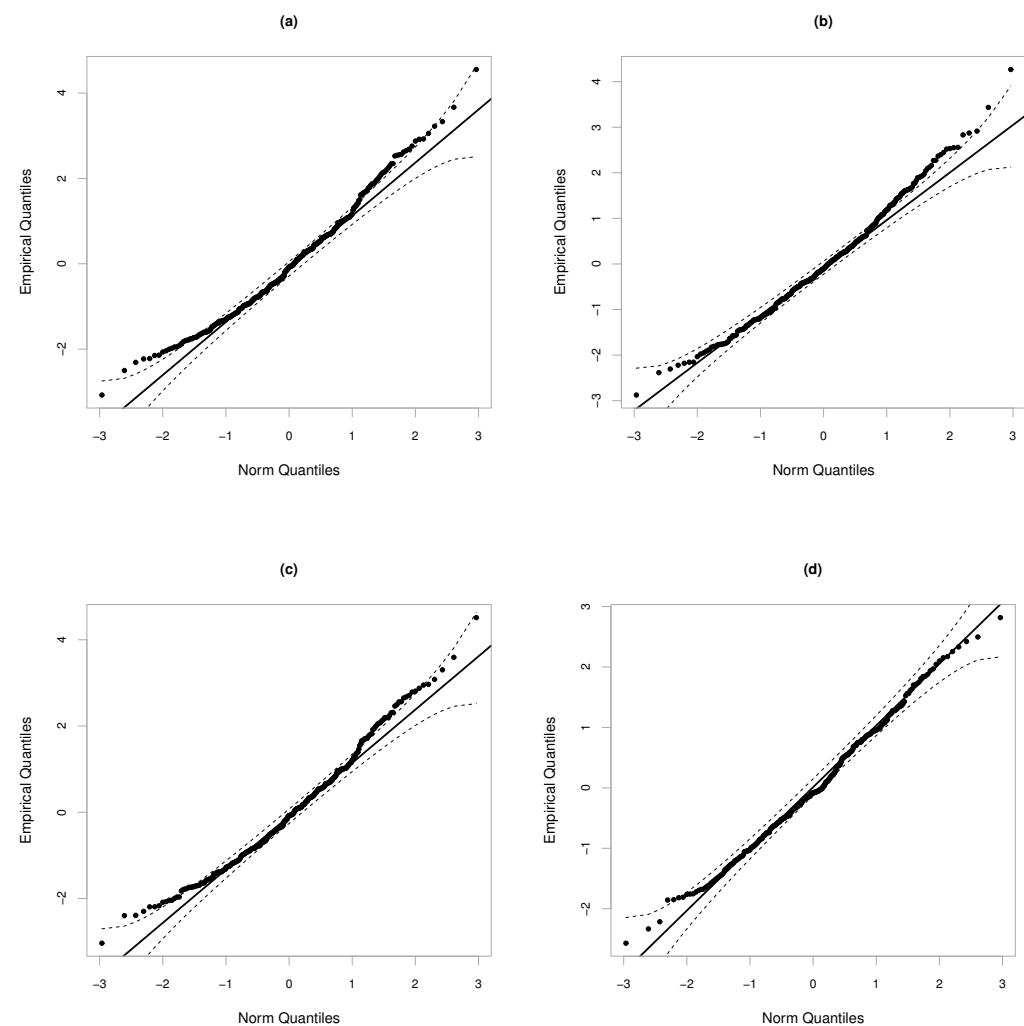
It is important to note that both the individual and interaction effects are statistically significant, as the corresponding *p*-values (not shown here) are less than 0.05. Additionally, the estimates of $\beta_0$ are similar across the four models, whereas the estimates of $\beta_1$ vary considerably, particularly in Model IV. Concerning the associated standard errors, all
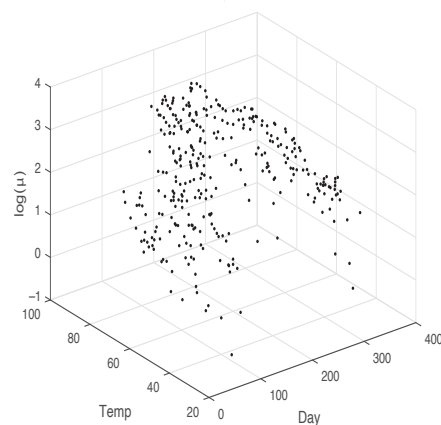
the estimators exhibit small values. The last two rows of Table 5 display the Akaike Information Criterion (AIC) and $R^2$ values, respectively. It is evident that the TPS-GLM, with $\text{AIC}(\lambda_f) = 1777.705$, provides the best fit to the Ozone data, followed by Model II with an AIC of 1806.837. This is corroborated by the QQ-plots in Figure 6, specifically Figure 6b,d. Furthermore, the $R^2$ value associated with our model is higher than those of Models I, II, and III. The smoothing parameter $\lambda_f$ was chosen such that the effective degrees of freedom were approximately 7. Figure 7 illustrates the 3D plot of the adjusted log(O3) against the explanatory variables Temp and Day, showing an adequate fit of the TPS-GLM.

**Table 5.** AIC, R$^2$(Adj), ML and MPL estimates for all four fitted models to the Ozone data.

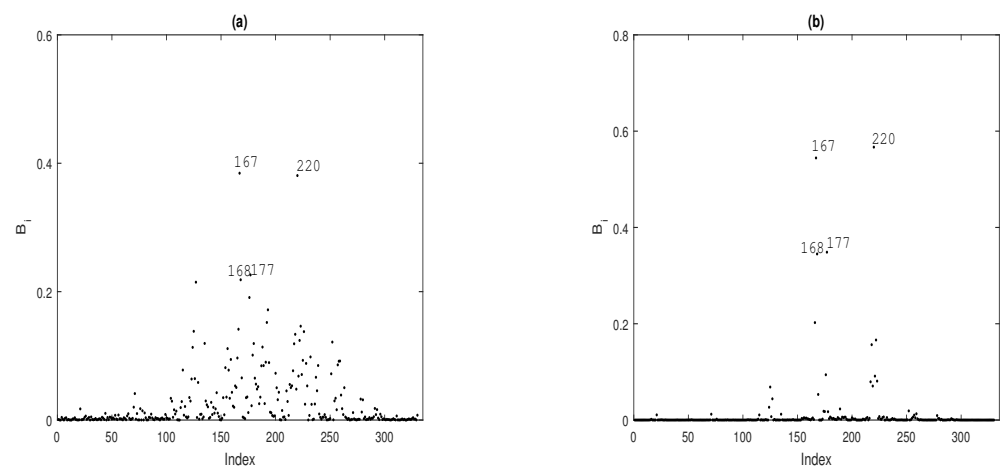| Parameters | I | II | III | IV |
|---|---|---|---|---|
| $\beta_0$ | 0.577 (0.104) | 0.478 (0.142) | 0.787 (0.198) | 2.507 (0.040) |
| $\beta_1$ | −0.002 (0.0003) | −0.002 (0.0003) | −0.002 (0.0003) | −0.002 (0.0003) |
| $\beta_2$ | 0.035 (0.001) | 0.033 (0.002) | 0.032 (0.003) | - |
| $\beta_3$ | −0.001 (0.002) | - | −0.002 (0.001) | - |
| $\beta_4$ | - | - | 0.00002 (0.00002) | - |
| AIC | 1887.312 | 1806.837 | 1887.757 | 1789.92 |
| R$^2$(Adj) | 0.673 | 0.715 | 0.670 | 0.728 |



**Figure 6.** QQ-plot of the standardized residuals for the models described in Table 5: Model I (**a**), Model II (**b**), Model III (**c**) and Model IV (**d**).

**Figure 7.** 3D plot between $\widehat{\log(\mu)}$ and explanatory variables Temp and Day.
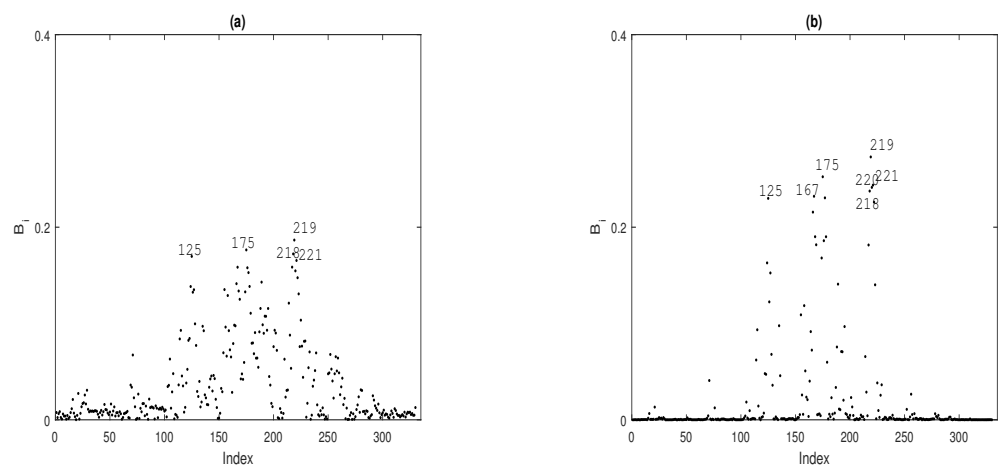
### 5.2.1. Diagnostic Analysis

To identify potentially influential observations on the MPL estimators under the fitted TPS-GLM for the Ozone data, we present some index plots of $B_i = B_{e_i}(\gamma)$, for $\gamma = \beta, \delta$. Figure 8 shows the index plot of $B_i$ for the case-weight perturbation scheme under the fitted model. In Figure 8a,b, note that observations #167, #220, #168, and #177 are more influential on $\widehat{\beta}$ and $\widehat{\delta}$, respectively. By perturbing the response variable additively, it becomes clear that observations #125, #175, #218, #219, and #221 are more influential on $\widehat{\beta}$ and $\widehat{\delta}$; see Figure 9a and 9b, respectively.



**Figure 8.** Index plots of $B_i$ for assessing local influence on $\widehat{\beta}$ (**a**) and $\widehat{\delta}$ (**b**), considering case-weight perturbation under model fitted to Ozone data.

From the local influence analysis, we conclude that the MPLE of the regression coefficients and the smooth surface are sensitive to perturbations in the data or the model. Furthermore, this analysis revealed that observations identified as influential for the parametric component are also influential for the nonparametric component, and vice versa. For example, under the case-weight perturbation scheme, observations #167, #220, #168, and #177 were found to be influential for both the parametric and nonparametric components.

**Figure 9.** Index plots of $B_i$ for assessing local influence on $\widehat{\boldsymbol{\beta}}$ (**a**) and $\widehat{\boldsymbol{\delta}}$ (**b**), considering response variable additive perturbation.

5.2.2. Confirmatory Analysis

To investigate the impact on model inference when influential potentially observations detected in the diagnostic analysis are removed, we present the relative changes (RCs) in the MPL estimate of $\beta_j$ for $j \in \{1, 2\}$ after removing the influential observations from the dataset (%). The RC is defined as $\text{RC}_\xi = \left| \frac{\widehat{\xi} - \widehat{\xi}_{(I)}}{\widehat{\xi}} \right| \times 100\%$, where $\widehat{\xi}_{(I)}$ denotes the MPL estimate of $\xi$, with $\xi = \beta_j$, after the corresponding observation(s) are removed according to set I. Table 6 presents the RCs in the regression coefficient estimates after removing the observations identified as potentially influential for the parametric component of the model.

**Table 6.** Relative changes (RCs) (in %) in the MPL estimates of $\beta_j$ under the TPS-GLM. The last two columns indicate the AIC and $R^2$(Adj) of the model with dropped observations.

| Dropped Obs. | $\beta_0$ | $\beta_1$ | $RC_{\beta_0}$ | $RC_{\beta_1}$ | AIC | $R^2$(Adj) |
|---|---|---|---|---|---|---|
| 167 | 2.513 | $-0.002$ | 0.231 | 0.378 | 1777.07 | 0.737 |
| 175 | 2.506 | $-0.002$ | 0.051 | 1.673 | 1784.58 | 0.724 |
| 219 | 2.540 | $-0.002$ | 1.334 | 7.105 | 1784.44 | 0.725 |
| 220 | 2.507 | $-0.002$ | 0.012 | 0.263 | 1777.87 | 0.728 |
| 167-175 | 2.511 | $-0.002$ | 0.169 | 1.052 | 1771.76 | 0.734 |
| 167-219 | 2.538 | $-0.002$ | 1.242 | 6.853 | 1771.58 | 0.735 |
| 167-220 | 2.511 | $-0.002$ | 0.179 | 0.884 | 1765.08 | 0.738 |
| 175-219 | 2.538 | $-0.002$ | 1.248 | 7.368 | 1779.10 | 0.722 |
| 175-220 | 2.504 | $-0.002$ | 0.104 | 0.684 | 1772.55 | 0.725 |
| 219-220 | 2.507 | $-0.002$ | 0.007 | 0.289 | 1772.807 | 0.725 |
| 167-175-219 | 2.536 | $-0.002$ | 1.155 | 7.136 | 1766.269 | 0.732 |
| 167-175-220 | 2.512 | $-0.002$ | 0.215 | 3.415 | 1759.79 | 0.735 |
| 175-219-220 | 2.504 | $-0.002$ | 0.098 | 0.678 | 1767.484 | 0.721 |
| 167-175-219-220 | 2.534 | $-0.002$ | 1.072 | 7.800 | 1754.761 | 0.731 |

## 6. Concluding Remarks and Future Research

In this work, we study some aspects of the Thin-Plate Spline Generalized Linear Models. Specifically, we derive an iterative process to estimate the parameters and the Fisher information matrix to approximate, through its inverse, the variance–covariance matrix of the estimators. In addition, we extended the local influence method, obtaining closed expressions for the Hessian and perturbation matrices under cases-weight perturbation and additive perturbation of the response variable. We performed a statistical data analysis with two real data sets of the agronomic and environmental area. The study showed the advantage of incorporating a smooth surface to model the joint effect of a pair of explanatory variables or the spatial effect determined by the coordinates. In both applications,

it was observed that the adjusted values of the response variable were consistent. In addition, it was observed that our model presented a better fit to model the soybean yield and ozone concentration data, compared to some classic parametric and semiparametric models, respectively. In our analysis, it was found that those observations detected as potentially influential generated important changes in the estimates, but not significant inferential changes. In addition, our study confirms the need to develop the local influence method to evaluate the sensitivity of maximum penalized likelihood estimators and thus determine those observations that can exert an excessive influence on both the parametric and non-parametric components, or on both.

As future work, we propose to incorporate a correlation component in the model and extend the local influence technique to other perturbation schemes, mainly on the non-parametric component of the model.

**Author Contributions:** Conceptualization, G.I.-P., P.P., M.A.U.-O. and O.N.; methodology, G.I.-P. and O.N.; software, G.I.-P. and P.P.; validation, G.I.-P. and P.P.; formal analysis, P.P.; investigation, G.I.-P., P.P., O.N. and M.A.U.-O.; data curation, P.P. and M.A.U.-O.; writing—original draft preparation, G.I.-P. and P.P.; writing—review and editing, O.N. and M.A.U.-O.; supervision, G.I.-P. and O.N.; project administration, O.N.; funding acquisition, O.N. All authors have read and agreed to the published version of the manuscript.

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Data Availability Statement:** Data cannot be shared openly but are available on request from authors.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## References

1. McCullagh, P.; Nelder, J.A. *Generalized Linear Models*, 2nd ed.; Chapman and Hall: London, UK, 1989.
2. Duchon, J. Interpolation des fonctions de deux variables suivant le principe de la flexion des plaques minces. *RAIRO Anal. Numér.* **1976**, *10*, 5–12. [CrossRef]
3. Duchon, J. Splines minimizing rotation-invariant semi-norms in Sobolev spaces. *Lect. Notes Math.* **1977**, *57*, 85–100.
4. Bookstein, F.L. Principal warps: Thin-plate splines and decomposition of deformations. *IEEE Trans. Pattern Anal. Mach. Intell.* **1989**, *11*, 567–585. [CrossRef]
5. Chen, C.; Li, Y.; Yan, C.; Dai, H.; Liu, G. A Thin Plate Spline-Based Feature-Preserving Method for Reducing Elevation Points Derived from LiDAR. *Remote Sens.* **2015**, *7*, 11344–11371. [CrossRef]
6. Wahba, G. *Spline Models for Observational Data*; SIAM: Philadelphia, PA, USA, 1990.
7. Green, P.J.; Silverman, B.W. *Nonparametric Regression and Generalized Linear Models*; Chapman and Hall: Boca Raton, FL, USA, 1994.
8. Wood, S.N. Thin plate regression splines. *J. R. Stat. Soc. Ser. B (Methodol.)* **2003**, *65*, 95–114. [CrossRef]
9. Moraga, M.S.; Ibacache-Pulgar, G.; Nicolis, O. On an elliptical thin-plate spline partially varying-coefficient model. *Chil. J. Stat.* **2021**, *12*, 205–228.
10. Cook, R.D. Assessment of Local Influence. *J. R. Stat. Soc. Ser. B (Methodol.)* **1986**, *48*, 133–169. [CrossRef]
11. Thomas, W.; Cook, R.D. Assessing influence on regression coefficients in generalized linear models. *Biometrika* **1989**, *76*, 741–749. [CrossRef]
12. Ouwens, M.N.M.; Tan, F.E.S.; Berger, M.P.F. Local influence to detect influential data structures for generalized linear mixed models. *Biometrics* **2001**, *57*, 1166–1172. [CrossRef]
13. Zhu, H.; Lee, S. Local influence for incomplete-data models. *J. R. Stat. Soc. Ser. B* **2001**, *63*, 111–126. [CrossRef]
14. Zhu, H.; Lee, S. Local influence for generalized linear mixed models. *Can. J. Stat.* **2003**, *31*, 293–309. [CrossRef]
15. Espinheira, P.L.; Ferrari, P.L.; Cribari-Neto, F. Influence diagnostics in beta regression. *Comput. Stat. Data Anal.* **2008**, *52*, 4417–4431. [CrossRef]
16. Rocha, A.; Simas, A. Influence diagnostics in a general class of beta regression models. *TEST* **2001**, *20*, 95–119. [CrossRef]
17. Ferrari, S.; Spinheira, P.; Cribari-Neto, F. Diagnostic tools in beta regression with varying dispersion. *Stat. Neerl.* **2011**, *65*, 337–351. [CrossRef]
18. Ferreira, C.S.; Paula, G.A. Estimation and diagnostic for skew-normal partially linear models. *J. Appl. Stat.* **2017**, *44*, 3033–3053. [CrossRef]
19. Emami, H. Local influence for Liu estimators in semiparametric linear models. *Stat. Pap.* **2018**, *59*, 529–544. [CrossRef]

20. Liu, Y.; Mao, G.; Leiva, V.; Liu, S.; Tapia, A. Diagnostic Analytics for an Autoregressive Model under the Skew-Normal Distribution. *Mathematics* **2020**, *8*, 693. [CrossRef]
21. Thomas, W. Influence diagnostics for the cross-validated smoothing parameter in spline smoothing. *J. Am. Stat. Assoc.* **1991**, *9*, 693–698. [CrossRef]
22. Ibacache, G.; Paula, G.A. Local Influence for student-t partially linear models. *Comput. Stat. Data Anal.* **2011**, *55*, 1462–1478. [CrossRef]
23. Ibacache-Pulgar, G.; Paula, G.A.; Galea, M. Influence diagnostics for elliptical semiparametric mixed models. *Stat. Model.* **2012**, *12*, 165–193. [CrossRef]
24. Ibacache, G.; Paula, G.A.; Cysneiros, F. Semiparametric additive models under symmetric distributions. *Test* **2013**, *22*, 103–121. [CrossRef]
25. Zhang, J.; Zhang, X.; Ma, H.; Zhiya, C. Local influence analysis of varying coefficient linear model. *J. Interdiscip. Math.* **2015**, *3*, 293–306. [CrossRef]
26. Ibacache-Pulgar, G.; Reyes, S. Local influence for elliptical partially varying coefficient model. *Stat. Model.* **2018**, *18*, 149–174. [CrossRef]
27. Ibacache-Pulgar, G.; Figueroa-Zuñiga, J.; Marchant, C. Semiparametric additive beta regression models: Inference and local influence diagnostics. *REVSTAT-Stat. J.* **2019**, *19*, 255–274.
28. Cavieres, J.; Ibacache-Pulgar, G.; Contreras-Reyes, J. Thin plate spline model under skew-normal random errors: Estimation and diagnostic analysis for spatial data. *J. Stat. Comput. Simul.* **2023**, *93*, 25–45. [CrossRef]
29. Jeldes, N.; Ibacache-Pulgar, G.; Marchant, C.; López-Gonzales, J.L. Modeling Air Pollution Using Partially Varying Coefficient Models with Heavy Tails. *Mathematics* **2022**, *10*, 3677. [CrossRef]
30. Saavedra-Nievas, J.C.; Nicolis, O.; Galea, M.; Ibacache-Pulgar, G. Influence diagnostics in Gaussian spatial—Temporal linear models with separable covariance. *Environ. Ecol. Stat.* **2023**, *30*, 131–155. [CrossRef]
31. Sánchez, L.; Ibacache-Pulgar, G.; Marchant, C.; Riquelme, M. Modeling Environmental Pollution Using Varying-Coefficients Quantile Regression Models under Log-Symmetric Distributions. *Axioms* **2023**, *12*, 976. [CrossRef]
32. Green, P.J. Penalized Likelihood for General Semi-Parametric Regression Models. *Int. Stat. Rev.* **1987**, *55*, 245–259. [CrossRef]
33. Nelder, J.A.; Wedderburn, R.W.M. Generalized Linear Models. *J. R. Stat. Soc. Ser. A (Gen.)* **1972**, *135*, 370–384. [CrossRef]
34. Wood, S.N. *Generalized Additive Models: An Introduction with R*, 2nd ed.; Chapman and Hall/CRC: Boca Raton, FL, USA, 2017.
35. Akaike, H. Information theory as an extension of the maximum likelihood principle. In *Proceedings of the Second International Symposium on Information Theory*; Petrov, B.N., Csaki, F., Eds.; Academiai Kiado: Budapest, Hungary, 1973.
36. Escobar, L.A.; Meeker, W.Q. Assessing Influence in Regression Analysis with Censored Data. *Biometrics* **1992**, *48*, 507–528. [CrossRef] [PubMed]
37. Billor, N.; Loynes, R.M. Local influence: A new approach. *Comm. Statist. Theory Meth.* **1993**, *22*, 1595–1611. [CrossRef]
38. MathWorks Inc. *MATLAB Version: 9.13.0 (R2022b)*; The MathWorks Inc.: Natick, MA, USA, 2022. Available online: https://www.mathworks.com (accessed on 10 October 2022).
39. Uribe-Opazo, M.A.; Borssoi, J.A.; Galea, M. Influence diagnostics in Gaussian spatial linear models. *J. Appl. Stat.* **2012**, *3*, 615–630. [CrossRef]
40. Breiman, L.; Friedman, J.H. Estimating optimal transformations for multiple regression and correlation. *J. Am. Stat. Assoc.* **1985**, *80*, 580–598. [CrossRef]
41. Ibacache-Pulgar, G.; Lira, V.; Villegas, C. Assessing Influence on Partially Varying-coefficient Generalized Linear Model. *REVSTAT-Stat. J.* **2022**. Available online: https://revstat.ine.pt/index.php/REVSTAT/article/view/507 (accessed on 10 October 2022).