On the Pulling Linear Regression and Its Applications in Digital Mammograms

NAHATAI TEPKASETKUL¹, WEENAKORN IEOSANURAK¹, THANAPONG INTHARAH², WATCHARIN KLONGDEE¹ ¹Department of Mathematics, Faculty of Science, Khon Kaen University, THAILAND ²Department of Statistics, Faculty of Science, Khon Kaen University, THAILAND

Abstract: - Regression analysis is a statistical approach used to investigate the correlations between variables, especially linear regression, that is a simple but effective approach for analyzing the relationship between a dependent variable and one independent variable. Since it has limitations based on the assumption that the mean of the noise should be zero, there are still some areas where it may be improved. In this article, we introduce a novel data fitting algorithm called the pulling linear regression, which is separated into two types: the line-pulling linear regression and the band-pulling linear regression. The method is developed from linear regression, which can create the regression line from the function that uses noise with various distributions. The result demonstrates that the sequence of sum square errors of the pulling linear regression is convergent. Moreover, we have a numerical example to show that the performance of the proposed algorithm is better than that of linear regression when the mean of the noise is not zero. And the last, we have an application to smooth the boundary of the pectoral muscle in digital mammograms. We found that the regression line of the proposed algorithm can do better than the linear regression when we would like to remove only the muscle part.

Key-Words: - linear regression, pulling linear regression, sum square error, root mean square error

Received: April 29, 2022. Revised: January 15, 2023. Accepted: February 9, 2023. Published: March 2, 2023.

1 Introduction

Least squares linear regression was performed by Legendre and Gauss for the prediction of planetary movement, [1]. Regression analysis is a statistical technique for examining correlations between variables used in numerous domains, including economics, engineering, physical science, biological science, social science, and medicine, among many others, [2].

Linear regression (LR) is a powerful and adaptable method for dealing with regression difficulties. The model definition, model estimation, statistical inference, model diagnosis, variable described selection, and prediction are comprehensively, [3]-[6]. Therefore, researchers are quite interested in the trend in LR models. For example, Pérez-Domínguez et al., [7], offered a contribution using linear regression and applied Dimensional Analysis (DA) to solve instability and error problems of the data transformation. Jokubaitis and Leipus, [8], studied the asymptotic normality in a high-dimensional linear regression where the covariance matrix of the regression variables has a KMS structure. Al-Kandari et al., [9], introduced a strategy for accounting for uncertainty in the

residuals of the linear regression model using fuzzy statistics. Liu and Chen, [10], improved the h value for fuzzy linear regression analysis using symmetric triangular fuzzy numbers and the least fuzziness criterion. Kabán, [11], provide a new analysis of compressive least squares regression that eliminates a false log N component, where N is the total number of training points. Additionally, several researchers have developed methods related to linear regression. For example, linear mixed models are used by Yi and Tang, [12], Ahn, Zhang and Lu, [13], and multiple linear regression is used by Uyanık and Güler, [14], Liu et al., [15], Li, He and Liu, [16].

A linear regression model is defined by $y = ax + b + \eta$ where *a* is a scaling parameter, *b* is a location parameter, *y* is the dependent or response variable, *x* is the independent or predictor variable, and the random variable η is the error term in the model, [17]-[19]. The linear regression is carried out under the assumption that η has a normal distribution with a mean and variance of zero and σ^2 , respectively, i.e., $E[\eta] = 0$ and $Var(\eta) = \sigma^2$.

In this article, we will consider the scenario where η has an alternative distribution or when

 $E[\eta] \neq 0$, which reduces the influence of partial observations by deviating from the regression line, y = ax + b. The two novel algorithms are line-pulling linear regression (LPR) and band-pulling linear regression (BPR), which are presented in the next section.

The remainder of the article is structured as follows: the LPR and BPR algorithms are introduced in Section 2. Section 3 describes a mathematical proof of some property. Next, the numerical results of our algorithms are illustrated and discussed in Section 4. Section 5 shows how the application is used to remove the pectoral muscle. Finally, conclusions and some suggestions are drawn and presented in Section 6.

2 Description of LPR and BPR

We introduce two novel data-fitting algorithms: line-pulling linear regression (LPR) and bandpulling linear regression (BPR). These algorithms are defined as follows:

Let $D^{(0)} = \{(x_1, y_1^{(0)}), (x_2, y_2^{(0)}), \dots, (x_n, y_n^{(0)})\}$ be an initial data such that x_i are distinct and $0 \le p_1, p_2 \le 1$. The procedure for the LPR algorithm is the following. Set $D_{LPR}^{(0)} = D^{(0)}$.

(i) Consider the k^{th} iteration, k = 1, 2, ... We get the linear regression $f^{(k)}(x) = \alpha^{(k)}x + \beta^{(k)}$ of the data $D_{LPR}^{(k-1)}$, where

$$(\alpha^{(k)}, \beta^{(k)}) = \underset{(\alpha, \beta)}{\operatorname{argmin}} \sum_{i=1}^{n} (y_i^{(k-1)} - (\alpha x_i + \beta))^2.$$

Denote $u_i^{(k)} = \alpha^{(k)} x_i + \beta^{(k)}, i = 1, 2, ..., n.$

- (ii) Define a band $M^{(k)} = \left\{ (x, y) \middle| -d_2^{(k)} \le y - f^{(k)}(x) \le d_1^{(k)} \right\},$ where $d_1^{(k)} = p_1 \max_i \left\{ y_i^{(k-1)} - u_i^{(k)} \right\},$ and $d_2^{(k)} = p_2 \max_i \left\{ u_i^{(k)} - y_i^{(k-1)} \right\}.$
- (iii) Update the data $D_{LPR}^{(k)}$ given by $D_{LPR}^{(k)} = \{ (x_1, y_1^{(k)}), (x_2, y_2^{(k)}), \dots, (x_n, y_n^{(k)}) \},$ such that for each $i = 1, 2, \dots, n,$ $(x_i, y_i^{(k)}) = \{ (x_i, y_i^{(k-1)}); (x_i, y_i^{(k-1)}) \in M^{(k)}, (x_i, y_i^{(k)}); (x_i, y_i^{(k-1)}) \notin M^{(k)}.$ (1)
- (iv) Return to step 1 to repeat until the error values

$$\sum_{i=1}^{n} \left(y_i^{(k-1)} - \left(\alpha^{(k)} x_i + \beta^{(k)} \right) \right)^2 < \varepsilon,$$

where $\varepsilon > 0$ is a fixed value.

Following that, we shall introduce the BPR algorithm. Set $D_{BPR}^{(0)} = D^{(0)}$. The steps of the algorithm are defined the same as the LPR algorithm, except for step (iii), replaced by step (iii*) as follows.

(iii*) Update the data
$$D_{BPR}^{(k)}$$
 given by
 $D_{BPR}^{(k)} = \left\{ \left(x_1, y_1^{(k)} \right), \left(x_2, y_2^{(k)} \right), \dots, \left(x_n, y_n^{(k)} \right) \right\},$
such that for each $i = 1, 2, \dots, n,$
 $\left(x_i, y_i^{(k)} + d_1^{(k)} \right); \quad (x_i, y_i^{(k-1)}) \notin M^{(k)},$
 $y_i^{(k-1)} > u_i^{(k)} + d_1^{(k)},$
 $\left(x_i, y_i^{(k-1)} \right); \quad (x_i, y_i^{(k-1)}) \in M^{(k)},$
 $\left(x_i, u_i^{(k)} - d_2^{(k)} \right); \quad (x_i, y_i^{(k-1)}) \notin M^{(k)},$
 $y_i^{(k-1)} < u_i^{(k)} - d_2^{(k)}.$
(2)

The two above algorithms are different for updating data. The LPR algorithm is used to pull the points outside the band toward the regression line, but the BPR algorithm is used to pull those points toward the boundary of the band as illustrated in Fig. 1.



Fig. 1: The proposed algorithm.

The following is an example of how to understand our algorithm.

Example 1. Let $p_1 = p_2 = 0.5$. We take $y_i = x_i + 1 + \eta_i$ where $x_i = i$ and $\eta_i \in (-1,1)$ is a uniform noise.

We obtain the initial data for our example, $D^{(0)} = \{(1,1.88), (2,2.76), (3,4.53), (4,5.59), (5,5.37)\},\$ (see, Fig. 2).

In the first iteration, we obtain the linear least square regression:

 $f^{(1)}(x) = 0.98x + 1.08.$ The points (1,2.06), (2,3.05), (3,4.03), (4,5.01), (5,5.99) that lie on the regression $y = f^{(1)}(x)$ and $d_1^{(1)} = 0.5 \max\{(1.88 - 2.06), (2.76 - 3.05), ($ (4.53 - 4.03), (5.59 - 5.01),(5.37 - 5.99) = 0.29
$$\begin{split} d_2^{(1)} &= 0.5\max\{(2.06-1.88), (3.05-2.76), \\ &(4.03-4.53), (5.01-5.59), \end{split}$$
(5.99 - 5.37) = 0.31 $M^{(1)} =$

Then,

 $\{(x, y) | -0.31 \le y - f^{(1)}(x) \le 0.29\}$ is shown in Fig. 3.



Fig. 2: Initial data.



Fig. 3: The band $M^{(1)}$ for the example.

In Fig. 3, we found that points $(x_3, y_3), (x_4, y_4)$ and (x_5, y_5) are outside of the band $M^{(1)}$. From equations (1) and (2), we will update the data as follows:

LPR: We get the range of $D_{LPR}^{(1)}$ as $\{y_1^{(0)}, y_2^{(0)},$ $u_3^{(1)}, u_4^{(1)}, u_5^{(1)}$ and obtain the updated data, $D_{LPR}^{(1)} = \{(1,1.88), (2,2.76), (3,4.03), (4,5.01), (5,5.99)\},\$ (see, Fig. 4(a)).

BPR: We get the range of $D_{BPR}^{(1)}$ as $\{y_1^{(0)}, y_2^{(0)},$

 $u_3^{(1)} + d_1^{(1)}, u_4^{(1)} + d_1^{(1)}, u_5^{(1)} - d_2^{(1)}$ and obtain the updated data, $D_{BPR}^{(1)} = \{(1,1.88), (2,2.76), (3,4.32), (4,5.30), (5,5.68)\},\$ (see, Fig. 5(a)).

The results of the LPR and BPR algorithms are shown in Tables 1 and 2, respectively. In the 6^{th} iteration for LPR and the 11th iteration for BPR, we obtain the regressions f(x) = 1.03x + 0.84 and f(x) = 1.04x + 0.84, respectively, for $\varepsilon = 0.5 \times$ 10^{-5} , shown in Fig. 4(b) and Fig. 5(b).



(a) the updated data for the 1^{st} iteration



(b) the data fitting for 6^{th} iteration

Fig. 4: The result for the LPR algorithm.



(a) the updated data for the 1^{st} iteration



(b) the data fitting for 11^{th} iteration

Fig. 5: The result for the BPR algorithm.

Table 1. The results of the LPR algorithm for $D^{(0)} = \{(1,1.88), (2,2.76), (3,4.53), (4,5.59), (5,5.37)\}.$

| k | $u_i^{(k)} = a^{(k)}$ | $x^{(k)}x_i + b^{(k)}$ | | Error | | | | |
|---|-----------------------|------------------------|-------------|-------------|---------------|-------------|---------------|------------------------|
| | $a^{(k)}$ | $b^{(k)}$ | $y_1^{(k)}$ | $y_2^{(k)}$ | $y_{3}^{(k)}$ | $y_4^{(k)}$ | $y_{5}^{(k)}$ | LIIOI |
| 0 | | | 1.880 | 2.760 | 4.530 | 5.590 | 5.370 | |
| 1 | 0.981 | 1.083 | 1.880 | 2.760 | 4.026 | 5.007 | 5.988 | 0.109×10 |
| 2 | 1.046 | 0.793 | 1.880 | 2.886 | 3.932 | 5.007 | 5.988 | 0.284×10^{-1} |
| 3 | 1.034 | 0.838 | 1.880 | 2.905 | 3.932 | 4.972 | 6.006 | 0.200×10^{-2} |
| 4 | 1.032 | 0.843 | 1.875 | 2.905 | 3.939 | 4.972 | 6.003 | 0.865×10^{-4} |
| 5 | 1.032 | 0.842 | 1.874 | 2.907 | 3.939 | 4.971 | 6.003 | 0.523×10^{-5} |
| 6 | 1.032 | 0.842 | | | | | | 0.174×10^{-6} |

Table 2. The results of the BPR algorithm for $D^{(0)} = \{(1,1.88), (2,2.76), (3,4.53), (4,5.59), (5,5.37)\}.$

| k | $u_i^{(k)} = a^{(k)} x_i + b^{(k)}$ | | | Error | | | | |
|----|-------------------------------------|-----------|-------------|-------------|---------------|-------------|---------------|------------------------|
| | $a^{(k)}$ | $b^{(k)}$ | $y_1^{(k)}$ | $y_2^{(k)}$ | $y_{3}^{(k)}$ | $y_4^{(k)}$ | $y_{5}^{(k)}$ | Enor |
| 0 | | | 1.880 | 2.760 | 4.530 | 5.590 | 5.370 | |
| 1 | 0.981 | 1.083 | 1.880 | 2.760 | 4.318 | 5.299 | 5.679 | 0.109×10 |
| 2 | 1.014 | 0.946 | 1.880 | 2.806 | 4.152 | 5.166 | 5.847 | 0.362×10^{0} |
| : | : | : | : | : | : | : | : | : |
| 10 | 1.398 | 0.840 | 1.880 | 2.919 | 3.961 | 5.000 | 6.038 | 0.128×10^{-4} |
| 11 | 1.398 | 0.840 | | | | | | 0.328×10^{-5} |

3 Main Results

Definition 1. Let $p_1, p_2 \in [0,1]$, $D^{(0)} = \{(x_1, y_1^{(0)}), (x_2, y_2^{(0)}), \dots, (x_n, y_n^{(0)})\}$ be an initial data such that x_i are distinct. The sum square error (SSE) of k^{th} iteration for LPR with respect to $(D^{(0)}; p_1, p_2)$ is given by

$$SSE_{LPR}^{(k)}(D^{(0)}; p_1, p_2) = \sum_{i=1}^n \left(y_i^{(k-1)} - \left(\alpha^{(k)} x_i + \beta^{(k)} \right) \right)^2,$$
(3)

where

$$\left(\alpha^{(k)},\beta^{(k)}\right) = \operatorname*{argmin}_{(\alpha,\beta)} \sum_{i=1}^{n} \left(y_i^{(k-1)} - (\alpha x_i + \beta)\right)^2,$$

and $y_i^{(k)}$ is an updated value as mentioned in LPR.

Similarly, the SSE of k^{th} iteration for BPR with

respect to
$$(D^{(0)}; p_1, p_2)$$
 is given by
 $SSE_{BPR}^{(k)}(D^{(0)}; p_1, p_2) = \sum_{i=1}^{n} (y_i^{(k-1)} - (\alpha^{(k)}x_i + \beta^{(k)}))^2,$
(4)

where

$$\left(\alpha^{(k)},\beta^{(k)}\right) = \underset{(\alpha,\beta)}{\operatorname{argmin}} \sum_{i=1}^{n} \left(y_i^{(k-1)} - (\alpha x_i + \beta)\right)^2,$$

and $y_i^{(k)}$ is an updated value as mentioned in BPR.

In specific case, $p_1 = p_2 = 0$. We observe that $d_1^{(1)} = d_2^{(1)} = 0$, i.e., $M^{(1)} = \{(x, y) | y = f^{(1)}(x)\}$. That is, each point in $D^{(0)}$ is pulled into the line $y = f^{(1)}(x)$. Therefore, $SSE_{LPR}^{(k)}(D^{(0)}; p_1, p_2) = 0$ and $SSE_{BPR}^{(k)}(D^{(0)}; p_1, p_2) = 0$, $k \ge 2$. In the opposite case, $p_1 = p_2 = 1$. We then get $d_1^{(k)} = \max_i \left\{ y_i^{(k-1)} - u_i^{(k)} \right\}$ and $d_2^{(k)} = \max_i \left\{ u_i^{(k)} - y_i^{(k-1)} \right\}$. It means that all points are in $M^{(k)}$, i.e., $D^{(k)} = D^{(0)}$. Therefore, $SSE_{LPR}^{(k)}(D^{(0)}; p_1, p_2) = SSE_{LPR}^{(1)}(D^{(0)}; p_1, p_2)$ and $SSE_{BPR}^{(k)}(D^{(0)}; p_1, p_2) = SSE_{BPR}^{(1)}(D^{(0)}; p_1, p_2)$, $k \ge 2$.

From Tables 1 and 2, we observe that the sum square errors are decreasing for both algorithms. This leads to the following property.

Lemma 1. Let $D^{(0)} = \{(x_1, y_1^{(0)}), (x_2, y_2^{(0)}), \dots, (x_n, y_n^{(0)})\}$ be an initial data such that x_i are distinct. If $p_1, p_2 \in [0,1]$, the sequences $SSE_{LPR}^{(k)}(D^{(0)}; p_1, p_2)$ and $SSE_{BPR}^{(k)}(D^{(0)}; p_1, p_2)$ are decreasing on k.

Proof. For the k^{th} iteration where k = 1, 2, 3, ...Let $u_i^{(k)} = \alpha^{(k)} x_i + \beta^{(k)}, i = 1, 2, ..., n$, such that, $\left(\alpha^{(k)}, \beta^{(k)}\right) = \underset{(\alpha,\beta)}{\operatorname{argmin}} \sum_{i=1}^n \left(y_i^{(k-1)} - (\alpha x_i + \beta)\right)^2$

Case *LPR*: From equation (1), we get

$$y_i^{(k)} = \begin{cases} y_i^{(k-1)}, & (x_i, y_i^{(k-1)}) \in M^{(k)}, \\ u_i^{(k)}, & (x_i, y_i^{(k-1)}) \notin M^{(k)}. \end{cases}$$

Thus,

$$\begin{pmatrix} y_i^{(k)} - u_i^{(k)} \end{pmatrix}^2 \\ = \begin{cases} \left(y_i^{(k-1)} - u_i^{(k)} \right)^2, & \left(x_i, y_i^{(k-1)} \right) \in M^{(k)}, \\ 0, & \left(x_i, y_i^{(k-1)} \right) \notin M^{(k)}. \\ \le \left(y_i^{(k-1)} - u_i^{(k)} \right)^2. \end{cases}$$

Next, consider equation (3),

$$SSE_{LPR}^{(k+1)}(D^{(0)}; p_1, p_2)$$

$$= \sum_{i=1}^n \left(y_i^{(k)} - \left(\alpha^{(k+1)} x_i + \beta^{(k+1)} \right) \right)^2$$

$$= \min_{(\alpha,\beta)} \sum_{i=1}^n \left(y_i^{(k)} - \left(\alpha x_i + \beta \right) \right)^2$$

$$\leq \sum_{i=1}^n \left(y_i^{(k)} - \left(\alpha^{(k)} x_i + \beta^{(k)} \right) \right)^2$$

$$= \sum_{i=1}^n \left(y_i^{(k)} - u_i^{(k)} \right)^2$$

$$\leq \sum_{i=1}^n \left(y_i^{(k-1)} - u_i^{(k)} \right)^2$$

$$= E_{LPR}^{(k)}(D^{(0)}; p_1, p_2).$$



$$y_i^{(k)} = \begin{cases} u_i^{(k)} + d_1^{(k)}, & \left(x_i, y_i^{(k-1)}\right) \notin M^{(k)}, \\ & y_i^{(k-1)} > u_i^{(k)} + d_1^{(k)}, \\ y_i^{(k-1)}, & \left(x_i, y_i^{(k-1)}\right) \in M^{(k)}, \\ u_i^{(k)} - d_2^{(k)}, & \left(x_i, y_i^{(k-1)}\right) \notin M^{(k)}, \\ & y_i^{(k-1)} < u_i^{(k)} - d_2^{(k)}. \end{cases}$$

Thus,

$$\begin{split} & \left(y_i^{(k)} - u_i^{(k)}\right)^2 \\ & = \begin{cases} \left(d_1^{(k)}\right)^2, & 0 \leq d_1^{(k)} < y_i^{(k-1)} - u_i^{(k)}, \\ \left(y_i^{(k-1)} - u_i^{(k)}\right)^2, & \left(x_i, y_i^{(k-1)}\right) \in M^{(k)}, \\ \left(-d_2^{(k)}\right)^2, & 0 \leq d_2^{(k)} < u_i^{(k)} - y_i^{(k-1)}, \\ \leq \left(y_i^{(k-1)} - u_i^{(k)}\right)^2. \end{split}$$

Next, consider equation (4),

$$SSE_{BPR}^{(k+1)}(D^{(0)}; p_1, p_2) = \sum_{i=1}^n \left(y_i^{(k)} - \left(\alpha^{(k+1)} x_i + \beta^{(k+1)} \right) \right)^2 \\ = \min_{(\alpha,\beta)} \sum_{i=1}^n \left(y_i^{(k)} - \left(\alpha x_i + \beta \right) \right)^2 \\ \le \sum_{i=1}^n \left(y_i^{(k)} - \left(\alpha^{(k)} x_i + \beta^{(k)} \right) \right)^2 \\ = \sum_{i=1}^n \left(y_i^{(k)} - u_i^{(k)} \right)^2 \\ \le \sum_{i=1}^n \left(y_i^{(k-1)} - u_i^{(k)} \right)^2 \\ \le E_{BPR}^{(k)}(D^{(0)}; p_1, p_2).$$

This completes the proof.

From equations (3) and (4), it is obvious that the sequences $SSE_{LPR}^{(k)}(D^{(0)}; p_1, p_2)$ and $SSE_{BPR}^{(k)}(D^{(0)}; p_1, p_2)$ are nonnegative sequences, that is, they have the lower bound to be zero, and, by Lemma 1, they are monotone decreasing, this leads to the following theorem.

Theorem 2. Let $D^{(0)}$ be an initial data. If $p_1, p_2 \in [0,1]$, the sequences $SSE_{LPR}^{(k)}(D^{(0)}; p_1, p_2)$ and $SSE_{BPR}^{(k)}(D^{(0)}; p_1, p_2)$ are convergent on k.

4 Numerical Examples

This section shows the numerical example using the initial data generated from the linear function y = x + 1 and noises (uniform distribution U(a, b), normal distribution N(a, b), and gamma distribution G(a, b) *). We also use various values of p_1 and p_2 to compare the root mean square error and the number of iterations of the proposed algorithm.

In Table 3, we generate y_i by the equation $y_i = (i - 11) + 1 + \eta_i$,

where η_i is a generated noise, and i = 1, 2, ..., 21. Thus, the root mean square error (RMSE) is given by

RMSE^(k) =
$$\sqrt{\frac{1}{21} \sum_{i=1}^{21} (u_i^{(k)} - y_i)^2}$$
,

where $u_i^{(k)}$ is obtained from the linear least square of k^{th} iteration of the proposed algorithms. We summarise the result in Table 3 as follows.

• $\eta_i \sim U(-10,0)$, the LPR algorithm with $p_1 = 1$ and $p_2 = 0.25$ has the minimum value of RMSE with 59th iteration.

- $\eta_i \sim U(-7.5,2.5)$, the BPR algorithm with $p_1 = 0.75$ and $p_2 = 0.25$ has the minimum value of RMSE with 15^{th} iteration.
- $\eta_i \sim U(-5,5)$, the LPR algorithm with $p_1 = 0.75$ and $p_2 = 0.75$ has the minimum value of RMSE with 11^{th} iteration.
- $\eta_i \sim U(-2.5,7.5)$, the BPR algorithm with $p_1 = 0.25$ and $p_2 = 0.75$ has the minimum value of RMSE with 17^{th} iteration.
- $\eta_i \sim U(0,10)$, the LPR algorithm with $p_1 = 0.25$ and $p_2 = 1$ has the minimum value of RMSE with 71th iteration.
- $\eta_i \sim N(-4,1)$, the LPR and BPR algorithms with $p_1 = 1$ and $p_2 = 0$ have the minimum value of RMSE with 52th iteration.

| | Algorithms | Proportions | | | | | | | | | | | |
|-------------|------------|---------------------|------------------------------|------------------------------|------------------------------|--|---------------------------|------------------------------|------------------------------|---------------------|---------------------------|------------------------------|------------------------------|
| Noises | | $p_1 = 0$ $p_2 = 0$ | $p_1 = 0.25$ $p_2 = 0.25$ | $p_1 = 0.50$ $p_2 = 0.50$ | $p_1 = 0.75$ $p_2 = 0.75$ | $\begin{array}{l} p_1=0\\ p_2=1 \end{array}$ | $p_1 = 0.25$ $p_2 = 1$ | $p_1 = 0.25$ $p_2 = 0.75$ | $p_1 = 0.50$ $p_2 = 0.75$ | $p_1 = 1$ $p_2 = 0$ | $p_1 = 1$ $p_2 = 0.25$ | $p_1 = 0.75$ $p_2 = 0.25$ | $p_1 = 0.75$ $p_2 = 0.50$ |
| | LPR | 3.555 | 3.623 | 3.627 | 3.663 | 8.414 | 8.414 | 4.259 | 3.966 | 0.144 | 0.143 | 2.860 | 3.198 |
| U(-10.0) | | (2) | (5) | (10) | (17) | (64) | (64) | (7) | (11) | (58) | (59) | (9) | (11) |
| 0(10,0) | BPR | 3.555 | 3.649 | 3.651 | 3.644 | 8.414 | 8.410 | 5.022 | 4.409 | 0.144 | 0.146 | 2.133 | 2.758 |
| | | (2) | (8) | (14) | (31) | (64) | (88) | (20) | (25) | (58) | (77) | (19) | (23) |
| | LPR | 1.563 | 1.562 | 1.509 | 1.269 | 5.819 | 5.819 | 2.202 | 1.676 | 2.357 | 2.357 | 0.401 | 0.889 |
| U(-7.5.2.5) | | (2) | (5) | (9) | (14) | (74) | (74) | (7) | (10) | (88) | (88) | (9) | (9) |
| - (-/ -/ | BPR | 1.563 | 1.511 | 1.464 | 1.425 | 5.819 | 5.815 | 3.266 | 2.606 | 2.357 | 2.330 | 0.174 | 0.337 |
| - | | (2) | (8) | (14) | (34) | (/4) | (93) | (16) | (22) | (88) | (122) | (15) | (19) |
| | LPR | 0.372 | 0.352 | 0.370 | 0.323 | 3.268 | 3.264 | 0.906 | 0.592 | 4.30/ | 4.36/ | 0.880 | 0.482 |
| U(-5,5) | | (2) | (5) | (7) | (11) | (05) | (07) | (8) | (12) | (88) | (88) | (10) | (11) |
| | BPR | (2) | (8) | (14) | (20) | 5.208 | (06) | (17) | (20) | 4.507 | 4.551 | 2.170 | (21) |
| | | (2) | (0) | (14) | (29) | 1.604 | (90) | 1 270 | (20) | 6 204 | 6 204 | 2.007 | (21) |
| | LPR | (2) | (5) | (8) | (15) | (55) | (55) | (9) | (11) | (63) | (64) | 2.007 | (11) |
| U(-2.5,7.5) | BPR | 1 698 | 1 721 | 1 701 | 1 673 | 1 694 | 1 692 | 0.203 | 0.718 | 6 304 | 6 299 | 3 346 | 2 729 |
| | | (2) | (8) | (13) | (29) | (55) | (72) | (17) | (23) | (63) | (87) | (18) | (21) |
| | LPR | 4.813 | 4.804 | 4.825 | 4.940 | 0.027 | 0.026 | 4.293 | 4.383 | 9.082 | 9.082 | 5,738 | 5.290 |
| | | (2) | (5) | (9) | (15) | (71) | (71) | (9) | (9) | (98) | (98) | (9) | (13) |
| U(0,10) | BPR | 4.813 | 4.803 | 4.839 | 4.914 | 0.027 | 0.075 | 3.110 | 3.765 | 9.082 | 9.067 | 6.477 | 5.890 |
| | | (2) | (8) | (14) | (32) | (71) | (103) | (19) | (21) | (98) | (132) | (18) | (21) |
| | IDD | 4.073 | 4.074 | 4.107 | 4.053 | 5.670 | 5.670 | 4.356 | 4.183 | 2.360 | 2.360 | 3.724 | 3.954 |
| N(-4.1) | LPK | (2) | (5) | (8) | (14) | (52) | (52) | (5) | (10) | (52) | (52) | (8) | (10) |
| N(-4,1) | BPR | 4.073 | 4.074 | 4.069 | 4.077 | 5.670 | 5.668 | 4.708 | 4.471 | 2.360 | 2.362 | 3.340 | 3.652 |
| | | (2) | (7) | (12) | (26) | (52) | (70) | (15) | (19) | (52) | (69) | (16) | (19) |
| | LPR | 1.235 | 1.186 | 1.058 | 1.040 | 2.459 | 2.460 | 1.295 | 1.207 | 0.877 | 0.874 | 0.863 | 0.903 |
| N(-11) | | (2) | (5) | (6) | (12) | (44) | (45) | (9) | (9) | (88) | (88) | (7) | (10) |
| 11(1,1) | BPR | 1.235 | 1.279 | 1.172 | 1.176 | 2.459 | 2.456 | 1.741 | 1.469 | 0.877 | 0.937 | 0.621 | 0.839 |
| | | (2) | (7) | (13) | (27) | (44) | (58) | (17) | (19) | (88) | (123) | (16) | (20) |
| | LPR | 1.082 | 1.078 | 1.108 | 1.079 | 1.200 | 1.201 | 0.702 | 0.915 | 3.013 | 3.014 | 1.396 | 1.317 |
| N(1.1) | | (2) | (5) | (8) | (13) | (57) | (58) | (9) | (10) | (53) | (54) | (7) | (9) |
| | BPR | 1.082 | 1.061 | 1.065 | 1.067 | 1.200 | 1.198 | 0.362 | 0.598 | 3.013 | 3.011 | 1.767 | 1.492 |
| | | (2) | (7) | (12) | (20) | (57) | (76) | (18) | (21) | (55) | (09) | (17) | (19) |
| | LPR | 3.899 | 3.900 | 3.927 | 4.095 | 2.007 | 2.009 | 3.640 | 5.795 (10) | 0.250 | 0.255 | 4.422 | 4.274 |
| N(4,1) | BPR | 3 800 | 3.046 | 3 072 | 3 000 | 2.007 | 2.014 | 3 185 | 3 520 | 6.256 | 6 253 | (7) | (12) |
| | | (2) | (7) | (13) | (27) | 2.007 | (87) | (17) | (22) | (57) | (79) | (19) | (21) |
| | | 0.788 | 0.788 | 0.788 | 0.739 | 0.023 | 0.022 | 0.637 | 0.693 | 1.662 | 1 663 | 0.873 | 0.828 |
| | LPR | (2) | (4) | (6) | (10) | (33) | (34) | (7) | (9) | (41) | (42) | (6) | (10) |
| G(1,1) | BPR | 0.788 | 0.782 | 0.776 | 0.777 | 0.023 | 0.025 | 0.389 | 0.528 | 1.662 | 1.659 | 1.175 | 1.022 |
| | | (2) | (7) | (13) | (26) | (33) | (44) | (13) | (19) | (41) | (54) | (13) | (19) |
| | LPR | 2.827 | 2.824 | 2.877 | 2.965 | 0.299 | 0.306 | 2.346 | 2.568 | 5.697 | 5.697 | 3.465 | 3.186 |
| | | (2) | (5) | (10) | (13) | (72) | (72) | (7) | (9) | (59) | (59) | (8) | (12) |
| 6(2,2) | DDD | 2.827 | 2.863 | 2.875 | 2.865 | 0.299 | 0.332 | 1.809 | 2.172 | 5.697 | 5.693 | 4.176 | 3.683 |
| | ВРК | (2) | (8) | (14) | (32) | (72) | (97) | (17) | (20) | (59) | (80) | (20) | (23) |

Table 3. The root mean square error and the number of iterations of the LPR and BPR.

* We random noises by using MATLAB program version R2017a, the codes for U(a, b), N(a, b), and G(a, b) are random('unif',a,b,m,n), random('norm',a,b,m,n), and random('gam',a,b,m,n), respectively, where $m \times n$ is a size of random noises.

- $\eta_i \sim N(-1,1)$, the BPR algorithm with $p_1 = 0.75$ and $p_2 = 0.25$ has the minimum value of RMSE with 16^{th} iteration.
- $\eta_i \sim N(1,1)$, the BPR algorithm with $p_1 = 0.25$ and $p_2 = 0.75$ has the minimum value of RMSE with 18^{th} iteration.
- $\eta_i \sim N(4,1)$, the LPR and BPR algorithms with $p_1 = 0$ and $p_2 = 1$ have the minimum value of RMSE with 65^{th} iteration.
- $\eta_i \sim G(1,1)$, the LPR algorithm with $p_1 = 0.25$ and $p_2 = 1$ has the minimum value of RMSE with 34^{th} iteration.
- $\eta_i \sim G(2,2)$, the LPR and BPR algorithms with $p_1 = 0$ and $p_2 = 1$ have the minimum value of RMSE with 72th iteration.

We observe that the number of iterations of the BPR algorithm is always greater than or equal to that of the LPR algorithm. Next, we discuss the values of p_1 and p_2 . In the case $p_1 = p_2$, the regression line located in the centre of the point. This designation is suitable for noise having a mean of zero. When $p_1 < p_2$, the bandwidth above the regression line is narrower than that below. The points that are pulled down are more than those that are pulled up. As a result, the regression line will gradually drop in the next iteration. This designation is suitable for noise with a mean of more than zero. When $p_1 > p_2$, the result is the opposite of the previous case. As a result, the regression line will gradually rise in the next iteration. This designation is appropriate for noise having a mean less than zero.

Furthermore, when every noise is positive, such as when η_i is derived from U(0,10), G(1,1), or G(2,2), the appropriate p is $p_2 = 1$. On the other hand, if every noise is negative, the appropriate p is $p_1 = 1$.

The following figures are some examples from Table 3. Fig. 6 shows the regression of the original function with $\eta_i \sim U(-5,5)$, $\eta_i \sim U(-2.5,7.5)$ and $\eta_i \sim U(-7.5,2.5)$, respectively. The black asterisks on the black line are the points (x, y) of the original function, the pink points are in $D^{(0)}$, and the red, green, and blue lines are the regression lines of linear regression, LPR, and BPR, respectively.



(a) $\eta_i \sim U(-5,5)$ and $p_1 = 0.75, p_2 = 0.75$



(b) $\eta_i \sim U(-2.5,7.5)$ and $p_1 = 0.25, p_2 = 0.75$



(c) $\eta_i \sim U(-7.5,2.5)$ and $p_1 = 0.75, p_2 = 0.25$ Fig. 6: The regression line with noise.

5 Applications

In general, the data points used to create the regression line have noise from the beginning. As a result, we cannot know whether those noises are positive or negative, making it impossible to choose a suitable *p*-value. Therefore, the selection of *p* must be determined depending on the desired outcome. For example, if the mean is to be used to represent the data, p_1 and p_2 should be equal. If the regression line is lower than the total data, $p_1 < 1$ and $p_2 = 1$ should be set. Similarly, set $p_1 = 1$ and $p_2 < 1$ when we need a regression line over all data points.

This section will demonstrate how the proposed algorithm can be used to analyse mammograms. In the mediolateral-oblique (MLO) view, the existence of the pectoral muscle may mislead the diagnosis of cancer due to its high-level similarity to the breast body. Therefore, we cut the pectoral muscle part and then employ LPR or BPR to smooth the boundary.

We separate the area of the pectoral muscle and breast using the difference in intensity of a mammogram. We have its border and transform to the Cartesian coordinate, which are referred to as the connection points and defined as $D^{(0)}$. Fig. 7 depicts an example of a mammogram and shows the connection points $D^{(0)}$ of the mammogram, with the area above the points representing the pectoral part and the area below representing the muscle part. Since we want to remove only the muscle part, the regression shall be below all points. We then define $p_2 = 1$.

The value of p_1 should be in the range [0,1); we specify $p_1 = 0.75$ for this example. Fig. 8 shows the regression lines derived from linear regression, LPR, and BPR using red, green, and blue, respectively. We get the regression y =-0.3383x + 464.8160 from linear regression, y =-0.4670 + 480.6891 from LPR in 2107th iteration, and y = -0.4752 + 484.2343 from BPR in 2705th iteration. Moreover, the mammograms after removing the pectoral muscle using linear regression, LPR, and BPR, respectively, are shown in Fig. 9.

The images obtained by finding the regression line by the LPR and BPR algorithms are comparable to and better than those obtained by the linear regression.



Fig. 7: An example of a mammogram.



Fig. 8: The regression lines.



(a) linear regression





(b) LPR

(c) BPR

Fig. 9: The mammograms after removing the muscle part by different algorithms.

6 Conclusions and Discussions

In this paper, we proposed the algorithm to create a regression line from an original function with noise η , where η is not necessarily a normal distribution with a mean of zero, called the line-pulling linear regression (LPR) and the band-pulling linear regression (BPR). These algorithms can set the regression line to the centre, top, or bottom of data points by assigning values p_1 and p_2 . If $p_1 = p_2 = 0$, the LPR and BPR algorithms provide the same regression lines as linear regression. When $p_1 < p_2$, the resulting line is below the linear regression. And when $p_1 > p_2$, the resulting line is above the linear regression. However, since we do not know the noise distribution in the data, we determine the value of p_1 and p_2 based on user requirements.

The numerical examples show that the results of the LPR and BPR algorithms are similar. The noticeable difference is the number of iterations for which the LPR algorithm converges faster than the BPR algorithm.

The application of these algorithms is the smoothing of the pectoral muscle's boundary. We use $p_1 = 0.75$ and $p_2 = 1$ to create the regression line at the bottom of all data points, ensuring we remove only the muscle part.

In addition, the LPR and BPR algorithms can be extended to more complicated models, such as using quadratic or cubic polynomial equations rather than a linear equation, which is expected to bring greater application benefits.

Acknowledgement:

The first author would like to express gratitude to the Science Achievement Scholarship of Thailand (SAST) for financial assistance for this paper. This research is supported by Department of Mathematics, Faculty of Science, Khon Kaen University, Fiscal Year 2022.

References:

- [1] S. M. Stigler, *The history of statistics: The measurement of uncertainty before 1900*, Harvard University Press, 1986.
- [2] D. C. Montgomery, E. A. Peck, G. G. Vining, *Introduction to linear regression analysis*, John Wiley & Sons, 2021.
- [3] X. Su, X. Yan, C. Tsai, Linear regression, *WIREs Computational Statistics*, Vol.4, No.3, 2012, pp. 275-294.
- [4] W. Yao, L. Li, A New Regression Model: Modal Linear Regression, *Scandinavian*

Journal of Statistics, Vol.41, 2014, pp. 656-671.

- [5] K. H. Zou, K. Tuncali, S. G. Silverman, Correlation and simple linear regression, *Radiology*, Vol.227, No.3, 2003, pp. 617-628.
- [6] D. Maulud, A. M. Abdulazeez, A review on linear regression comprehensive in machine learning, *Journal of Applied Science and Technology Trends*, Vol.1, No.4, 2020, pp.140-147.
- [7] L. Pérez-Domínguez, H. Garg, D. Luviano-Cruz, J.L. García Alcaraz, Estimation of Linear Regression with the Dimensional Analysis Method, *Mathematics*, Vol.10, No.10, 2022, pp. 1645.
- [8] S. Jokubaitis, R. Leipus, Asymptotic normality in linear regression with approximately sparse structure, *Mathematics*, Vol.10, No.10, 2022, pp. 1657.
- [9] M. Al-Kandari, K. Adjenughwure, K. Papadopoulos, A Fuzzy-Statistical Tolerance Interval from Residuals of Crisp Linear Regression Models, *Mathematics*, Vol.8, No.9, 2020, pp. 1422.
- [10] X. Liu, Y. Chen, A systematic approach to optimizing *h* value for fuzzy linear regression with symmetric triangular fuzzy numbers, *Mathematical Problems in Engineering*, Vol.2013, 2013.
- [11] A. Kabán, New bounds on compressive linear least squares regression, *Artificial intelligence and statistics*, 2014, pp. 448-456.
- [12] J. Yi, N. Tang, Variational Bayesian inference in high-dimensional linear mixed models, *Mathematics*, Vol.10, No.3, 2022, pp. 463.
- [13] M. Ahn, H. H. Zhang, W. Lu, Moment-based method for random effects selection in linear mixed models, *Statistica Sinica*, Vol.22, No.4, 2012, pp. 1539.
- [14] G. K. Uyanik, N. Güler, A Study on Multiple Linear Regression Analysis, *Procedia - Social* and Behavioral Sciences, Vol.106, 2013, pp. 234-240.
- [15] M. Liu, S. Hu, Y. Ge, G. B. Heuvelink, Z. Ren, X. Huang, Using multiple linear regression and random forests to identify spatial poverty determinants in rural China, *Spatial Statistics*, Vol.42, 2021, pp. 100461.
- [16] Y. Li, X. He, X. Liu, Fuzzy multiple linear least squares regression analysis, *Fuzzy Sets and Systems*, 2022.
- [17] S. Weisberg, *Applied Linear Regression*, 4th editio, 2014.

- [18] M. S. Paolella, *Linear models and time-series* analysis: regression, ANOVA, ARMA and GARCH, John Wiley & Sons, 2018.
- [19] A. C. Rencher, G.B. Schaalje, *Linear models in statistics*, John Wiley & Sons, 2008.

Contribution of Individual Authors to the Creation of a Scientific Article (Ghostwriting Policy)

-Nahatai Tepkasetkul carried out the conceptualization, investigation, methodology, software, validation, visualization, writing - original draft, and writing - review & editing.

-Weenakorn Ieosanurak carried out the supervision, validation, and writing - review & editing.

-Thanapong Intharah carried out the conceptualization, supervision, validation, and writing - review & editing.

-Watcharin Klongdee carried out the conceptualization, investigation, methodology, supervision, validation, visualization, writing - original draft, and writing - review & editing.

Sources of Funding for Research Presented in a Scientific Article or Scientific Article Itself

The first author would like to express gratitude to the Science Achievement Scholarship of Thailand (SAST) for financial assistance for this paper. This research is supported by Department of Mathematics, Faculty of Science, Khon Kaen University, Fiscal Year 2022.

Conflict of Interest

The authors have no conflicts of interest to declare that are relevant to the content of this article.

Creative Commons Attribution License 4.0 (Attribution 4.0 International, CC BY 4.0)

This article is published under the terms of the Creative Commons Attribution License 4.0 <u>https://creativecommons.org/licenses/by/4.0/deed.en</u> <u>US</u>