# NONPARAMETRIC ADAPTIVE SMOOTHING WITH PRACTICAL APPLICATION

**Noor Salah Hassan[1, a)]**
**Sahera Hussein Zain Al-Thalabi[2,b)]**
[1,2]*Department of Statistics, College of Administration and Economics, Basra University, Basra, Iraq*
a) *Corresponding author: noursalahh1997@gmail.com*
b) *sahera.zain@uobasrah.edu.iq*

## INTRODUCTION

There are many problems that the researcher faces when estimating the nonparametric regression functions because the estimation methods depend on the data, as these estimates may be inaccurate, or they may not be suitable for the nonparametric model, so the aim of this study is to find the adaptive capabilities in the nonparametric regression using modern methods. To increase the efficiency of estimation through the use of adaptive estimators in nonparametric regression smoothing . We will discuss some studies that used the adaptive method and its use in nonparametric regression, including:

-The study (Hill and others, 1988) used two nonparametric adaptive procedures to apply multiple comparisons and a test of alternatives required in a one-way ANOVA model, in comparison with the parametric normal theoretical procedure, and the rank-based non-parametric procedure where these procedures are applied to lung cancer data. The results showed the superiority of the adaptive procedures Nonparametric. [6]

-A study (2021, page and Grunewalder) presented an Adaptive estimation using the modern Goldenshluger-Lepski method to choose parameters for the statistical estimator using only the available data without making strong assumptions about the estimation. Nucleus . This method was used to address two regression problems, the kernel regression was fixed in one of them and in the other an adaptation was used. [12]

-The study (Breunig and Chen, 2022) aimed to find an adaptive estimation of the minimum quadratic function in the model of non-parametric automatic variables (NPIV), which is an important problem in the optimal estimation of non-linear functions, this problem is solved through a choice based on data from Lepski type For the smoothing parameter, the results showed that the adaptive estimator of the quadratic function achieves the minimum optimum rate. [3].Adaptive estimator in nonparametric regression: [10][1][3][4][8][11] An adaptive estimator is defined as an effective estimator for only a partially specified model ("effective" meaning that it is asymptotically equivalent to a non-parametric "likelihood Maximum" local probability estimator Applicable), or a model whose distribution is unknown, so adaptive estimation aims to build estimations entirely based on data without making strong assumptions about the estimation. Nonparametric regression is also a form of regression analysis and a common and flexible tool for data analysis and modeling of the non-linear relationship between dependent and explanatory variables. , that is, it depends mainly on the data, Where the objective of the nonparametric regression is to estimate the regression function without dependence or having prior knowledge of its functional form , and using adaptive methods, Classical methods can also be modified to be as robust as non-parametric methods . Studies to build a method for selecting data-based smoothing parameters in order to obtain adaptive estimates. The first adaptive estimate was proposed by (lepski 1990) and was developed in (1992) and its goal was to build capabilities from the data in the best possible way and reduce the risk of estimation, the adaptive methods in regression The non-parametric is strong in efficiency

as it cannot be outweighed by any non-adaptive method, as the exact adaptive procedure will work well with the data. So the adaptive approach is mainly divided into two types, The adaptive procedure for estimating unknown parameters is such as in a nonparametric regression, or the use of Data to determine the appropriate statistical procedure, the adaptive non-parametric approach on the one hand is estimating the parameters from the sample, or data-driven methods may be the best and most, The first to suggest this approach (Randles and Hogg). So the main purpose of adaptive approaches may be to provide a relatively easy alternative to parameterization without much effort on how to choose one from a variety of methods, and to facilitate the decision on the use of the appropriate technique. Adaptive approaches can perform better based on the available information in terms of achieving the desired combination of robustness and efficiency. over the past ten years. Adaptive-order tests show that adaptive actions Adaptive method can increase the power of tests, If the distribution of random error is abnormal, the power of classical tests is much lower than adaptive tests.

The formula for nonparametric regression is as follows:
$$, \quad i = 1,2,\ldots,n \ , \quad \varepsilon \sim N(0,\sigma^2) y_i = m(x_i) + \varepsilon_i \quad \ldots \ (1)$$
$Y_i$: the response variable, $m(x_i)$: the unknown function to be estimated, $x_i$: the explanatory variable , $\varepsilon_i$: the values of the random variable, which is white noise that is normally distributed.
The adaptive estimator for the parameter vector is as follows [15] :
$$\hat{\theta}(x) = \hat{\theta}_k(x) = \left( \theta_k^1(x), \ldots \ldots, \theta_k^p(x) \right)^T \ \ldots (2) \qquad k = 1, \ldots., p$$
$\theta_k^1, \ldots \ldots, \theta_k^p$ :Unknown parameter, θ is estimated based on sample observations $(x_i, y_i)$

## KERNEL SMOOTHERS
The positional polynomial regression smoother (LLS) is one of the best smoothing methods because it deals with static and random models, and it is sometimes called the weight or window function, as this function is continuous and symmetric, its integral is equal to the integer one, when (the bandwidth) is small very . [10]
The formula for smoothers is as follows
$$\widehat{m_h}(x) = \frac{\sum_{i=1}^{n} y_i \, k \, (x - X_i)/h}{\sum_{i=1}^{n} k \, (x - X_i)/h} \ \ldots (3)$$
$$w_i(x) = \frac{\dfrac{k(x - X_i)}{h}}{\sum_{i=1}^{n} k \, (x - X_i)/h} \ \ldots.. (4)$$

$\frac{\sum_{i=1}^{n} k(x-x_i)}{h}$ : represents the endodontic function , $w_i(x)$: represents the weight function and one of its conditions is positive, h: represents the smoothing parameter (the bandwidth) in the estimator $(m(x))$. If its value is large, the function is smooth, and if its value is small, the function is not smooth. [1]
-The Gasser-Müller (GM) smoother is one of the most widely used gradient smoothing tools. The Gasser-Müller estimator which is a modification of the Priestley-chao, estimator is used to construct nonparametric estimates of the regression function,), a new type of kernel. [4]
Its general form is as follows. [4]
$$\widehat{m_h}(x) = \frac{1}{n} \sum_{i=1}^{n} Y_i \int_{s_{i-1}}^{s_i} K\left(\frac{x-u}{h}\right) du \quad \ldots \quad (5) \quad s_0 = 0 \ , \quad s_n = 1$$
$$s_i = \frac{x_i + x_{i+1}}{2} \qquad , x_i \le s_i \le x_{i+1}$$

- Also, the nearest neighbor smoother (K-NN) depends on calculating the Euclidean distance between each point and the point closest to it. If the data are close to each other, the distance will be small and vice versa. [9]
So, its general form is as follows:
$$\widehat{m_k}(x) = \frac{\dfrac{k(x_i - x)}{k_l}}{\sum_{i=1}^{n} \dfrac{k(x_i - x)}{k_l}} \qquad k_l \to \infty \ldots (6)$$
$$, k_l = d(i,j) = \sqrt[2]{\sum_{i=1}^{n} (x_i - x_j)^2} \ \ldots \ (7)$$
$k_l$ :represents the Euclidean distance between x, k and : $x_i$ ,$x_j$ data points

## SPLINE SMOOTHER
depend on the sum of the squares of the error as used when the regression line is divided into pieces, as the explanatory variable x with period (a,b) is divided and the lines cut are called slide nodes so that smoothing the slides overcomes the problem of choosing a node and from During the identification of new nodes or changing the existing nodes, they are divided into linear spline (SPL) and cubic spline (SPC). [2][9]

$$S(m) = \sum_{i=1}^{n}\left(yi - \widehat{m}(xi)\right)^2 + \lambda \int_{a}^{b}[\widehat{m}^{`}(x)]^2\, dx \; \dots \; (8) \;, \quad \lambda > 0$$

Whereas

$\sum_{i=1}^{n}\left(yi - \widehat{m}(xi)\right)^2$: It represents the sum of the squares of the error ,
$\widehat{m}^{`}(x)$: Represents the second derivative of the bootstrap function , $\lambda$ : Represents the penalty factor indicating the width of the appropriateness quality package represented by $\sum_{i=1}^{n}\left(yi - \widehat{m}(xi)\right)^2$ And the smoothing of appreciation represented by $\int_{a}^{b}[\widehat{m}^{`}(x)]^2\, dx$

Goldenshluger-Lepski adaptive bandwidth extends Lepski's method for performing adaptation across multiple parameters .This method has been used in different contexts as it was used for the first time in a multidimensional white noise model. As it has been widely used in recent studies of non-parametric estimation, the idea of this method for adaptive non-parametric estimation is to choose an estimator that reduces the sum of the unknown bias factor of variance. [8][12]

The Goldenshluger-lepski formula is as follows [5] :

$$\hat{h}(x_i) = \arg min_{h \in H_n}\left\{\hat{A}(h, x_i) + \hat{V}(h, x_i)\right\} \; \dots \; (9)$$

$$\hat{A}(h, x_i) = max_{h' \in H_n}([\widehat{m}_{h'}(x_i) - \widehat{m}_{h \vee h'}(x_i)]^2 - V(h', x_i)) \dots (10)$$

$$\hat{V}(h, x_i) = k\sigma^2 \frac{\ln n}{n\widehat{\varphi}(h)} \; , h \neq 0 \dots (11)$$

K : represent a constant that does not depend on h , $\widehat{m}_h(x_i)$: function estimator , $H_n$ : Represents a set of smoothing parameter (bandwidth) .

$\hat{V}(h, x_i)$: Represents an empirical analogue of variance , $\hat{A}(h, x_i)$ : Represents an approximation of the term bias .

In order to estimate the regression curve, there are several criteria that are relied upon in the differentiation, and among these criteria are the mean absolute error squares (MAS), the roots mean squares error (RMSE), and the mean squared error (MSE) standard [9][14]. The function was used Endodontic (Epanchnickov) and adaptive bandwidth (Goldenshluger-lepski) on the experimental side.

$$\dots \qquad (12) \text{MAE} = \frac{1}{n}\sum_{i=1}^{n}|y_i - \widehat{m}(x)|$$

$$\text{RMSE} = \sqrt[2]{\frac{1}{n}\sum_{i=1}^{n}\left(y_i - \widehat{m}(x)\right)^2} \quad \dots \; (13)$$

$$MSE = \frac{1}{n}\sum_{i=1}^{n}(y_i - \widehat{m}(x))^2 \dots \; (14)$$

## STATISTICAL ANALYSIS

The analysis of this study is carried out using simulation, as it is known as a method that includes the use of a theoretical mathematical model and similar to the real model that represents the studied problem.  Simulation experiments were carried out using three sample sizes (n = 30, 60, 100) and with a frequency of 500 for each experiment. The nonparametric methods will be compared, and two models were used in the simulation.

first model $\qquad m(x_i) = 1 + 0.8e^{-200.(-0.5+x)^2} + 2x^2$

The variables (independent and random error) were generated. The random errors are normally distributed with a mean of zero and variance $\sigma^2$, The nonparametric explanatory variable Xi is generated according to the standard normal distribution .

$$e_i \sim N(0, \sigma^2)$$
$$X_i \sim N(0,1)$$

**TABLE 1**. The first model, (RMSE, MSE, MAE) criteria for the first model according to the different sample sizes and levels of variation

| $\sigma^2$ | n | ALLS | AGM | KNN | ASPL | ASPC |
|---|---|---|---|---|---|---|
| **RMSE** | | | | | | |
| $\sigma^2 = 0.5$ | 30 | 0.516265 | 0.930748 | 0.893772 | 0.928897 | 0.932278 |
| | 60 | 0.558081 | 0.861925 | 0.809613 | 0.842828 | 0.868399 |
| | 100 | 0.51343 | 0.759951 | 0.803071 | 0.833716 | 0.818787 |
| $\sigma^2 = 1$ | 30 | 1.304713 | 1.293053 | 1.391125 | 1.297777 | 1.300542 |
| | 60 | 1.033427 | 1.288025 | 1.256423 | 1.277987 | 1.282404 |
| | 100 | 1.134386 | 1.143293 | 1.373155 | 1.14562 | 1.155014 |
| $\sigma^2 = 1.5$ | 30 | 1.757319 | 1.737183 | 2.094056 | 1.742326 | 1.738227 |
| | 60 | 1.661094 | 1.739168 | 2.06102 | 1.716176 | 1.72202 |
| | 100 | 1.540448 | 1.666366 | 1.550798 | 1.664799 | 1.680947 |
| **MSE** | | | | | | |
| $\sigma^2 = 0.5$ | 30 | 0.266529 | 0.866292 | 0.798829 | 0.862849 | 0.869143 |
| | 60 | 0.311454 | 0.742915 | 0.655473 | 0.710359 | 0.754118 |
| | 100 | 0.263611 | 0.577525 | 0.610768 | 0.695083 | 0.670413 |
| $\sigma^2 = 1$ | 30 | 1.702275 | 1.671985 | 1.93523 | 1.684224 | 1.69141 |
| | 60 | 1.046111 | 1.659008 | 1.067971 | 1.633252 | 1.644561 |
| | 100 | 1.309619 | 1.307119 | 1.885555 | 1.312445 | 1.334058 |
| $\sigma^2 = 1.5$ | 30 | 3.088172 | 3.024758 | 4.385072 | 3.035701 | 3.021432 |
| | 60 | 2.759233 | 3.024707 | 4.247802 | 2.945262 | 2.965353 |
| | 100 | 2.69107 | 2.776777 | 2.404973 | 2.771554 | 2.825582 |
| **MAE** | | | | | | |
| $\sigma^2 = 0.5$ | 30 | 0.412093 | 0.718155 | 0.715842 | 0.730157 | 0.727187 |
| | 60 | 0.471807 | 0.651849 | 0.627654 | 0.639025 | 0.652287 |
| | 100 | 0.39879 | 0.63034 | 0.601175 | 0.68473 | 0.673663 |
| $\sigma^2 = 1$ | 30 | 0.997865 | 1.001918 | 1.116884 | 1.002198 | 1.014442 |
| | 60 | 0.955424 | 0.963265 | 1.013914 | 0.996628 | 0.993963 |
| | 100 | 0.851318 | 0.850425 | 0.826585 | 0.849166 | 0.867145 |
| $\sigma^2 = 1.5$ | 30 | 1.343314 | 1.409301 | 1.678374 | 1.380658 | 1.379656 |
| | 60 | 1.281802 | 1.288404 | 1.525217 | 1.289815 | 1.332256 |
| | 100 | 1.258361 | 1.264312 | 1.241343 | 1.288155 | 1.316453 |

Source/ From the (R.4.1.2) Package using simulation method

Explanation of Table 1. for the first model
1-The results showed, depending on the comparison criteria (RMSE) and (MSE) when the sample size is (n = 30, 60, 100) and with a level of variance ($\sigma^2 = 0.5$) that the best adaptive estimator is (ALLS), but when the level of variance is ($\sigma^2 = 1$ ,1.5) and sample size (n = 30), then the best adaptive estimator is ((AGM). As for (n=60 , $\sigma^2 = 1.5$) the best adaptive estimator is (ALLS), then the adaptive estimator (ASPL) .
2-The results showed that, depending on the comparison standard (MAE), when the sample size is (n=30,60,100) and with the level of variance $\sigma^2 = 0.5$, the estimator is (ALLS), then the estimator is (KNN), but when the variance level is ($\sigma^2 = 1$) At the sample size (n=30,60), the best adaptive estimator is (ALLS), followed by the adaptive estimator (AGM), and when the level of variance is ($\sigma^2 = 1.5$) at the sample size (n=30), the best estimator It is an adaptive estimator (ALLS), followed by an adaptive estimator (ASPC).
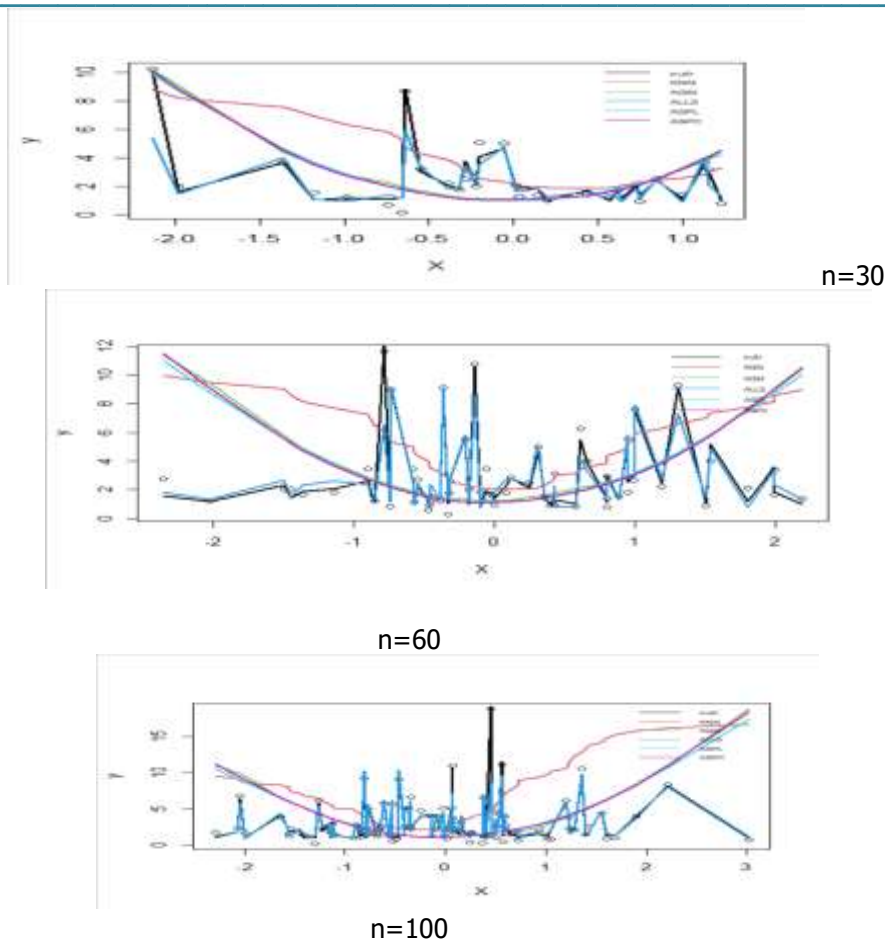
n=30



n=60



n=100

**FIGURE1.** The adaptive nonparametric capabilities of the first model when the sample size is (30, 60, 100)

## APPLIED ANALYSIS

After obtaining the best adaptive estimator (ALLS) for the nonparametric regression model on the experimental side (simulation) and for three different criteria, in this aspect the best adaptive estimator (ALLS) will be applied to the real data through the use of five explanatory variables and a response variable, The data was obtained from the Directorate of Civil Defense in Basra Governorate for a period of three years for the period (2018-2020).  And the study variables are as follows :   X1: represents the electrical contact    , X2: Represents a cigarette butt , X3: Represents children's tampering , The dependent variable (Y) represents the total number of fires .

**TABLE 2**  Estimated nonparametric adaptive functions

| Obs | $\widehat{m}(x_1)$ | $\widehat{m}(x_2)$ | $\widehat{m}(x_3)$ |
|---|---|---|---|
| 1 | 0.228887309 | 0.120492996 | 0.207268802 |
| 2 | 0.192097027 | 0.28166734 | 0.242724737 |
| 3 | 0.197002368 | 0.265104151 | 0.20702813 |
| 4 | 0.170260512 | 0.50441889 | 0.361497975 |
| 5 | 0.299637381 | 0.28166734 | 0.3057311 |
| 6 | 0.619300297 | 0.53478687 | 0.68408641 |
| 7 | 0.647639433 | 0.53478687 | 0.70582833 |
| 8 | 0.576975284 | 0.838661981 | 0.805051443 |
| 9 | 0.541237077 | 0.263131001 | 0.615711808 |
| 10 | 0.36774093 | 0.329188332 | 0.365872915 |
| 11 | 0.138214313 | 0.077805087 | 0.106955038 |
| 12 | 0.010981869 | 0.162130098 | 0.039099197 |
| 13 | 0.210813314 | 0.120492996 | 0.161488039 |
| 14 | 0.123128238 | 0.198472028 | 0.063214676 |
| 15 | 0.119379087 | 0.263131001 | 0.2217636 |
| 16 | 0.19829425 | 0.265104151 | 0.2217636 |
| 17 | 0.281285794 | 0.28166734 | 0.290890761 |

| | | | |
|---|---|---|---|
| 18 | 0.697523475 | 0.492733375 | 0.290890761 |
| 19 | 0.62397061 | 0.53478687 | 0.716552551 |
| 20 | 0.6219929 | 0.50441889 | 0.571542824 |
| 21 | 0.36774093 | 0.266575494 | 0.395207857 |
| 22 | 0.271036442 | 0.240113407 | 0.361246217 |
| 23 | 0.311749537 | 0.265104151 | 0.212019707 |
| 24 | 0.210813314 | 0.198472028 | 0.126618638 |
| 25 | 0.26318656 | 0.162130098 | 0.2217636 |
| 26 | 0.281285794 | 0.265104151 | 0.212019707 |
| 27 | 0.181764214 | 0.265104151 | 0.20702813 |
| 28 | 0.088569838 | 0.411743182 | 0.216996127 |
| 29 | 0.251460476 | 0.557401107 | 0.379633548 |
| 30 | 0.558298159 | 0.521875542 | 0.523286152 |
| 31 | 0.630571844 | 0.570845448 | 0.437083931 |
| 32 | 0.619300297 | 0.496713116 | 0.361246217 |
| 33 | 0.228887309 | 0.492733375 | 0.316031655 |
| 34 | 0.194355126 | 0.492733375 | 0.216996127 |
| 35 | 0.189988621 | 0.265104151 | 0.205834566 |
| 36 | 0.219517763 | 0.162130098 | 0.126618638 |

The criteria and coefficient of determination for each estimated function.

| $\widehat{m}(x)$ | $R^2$ | MAE | RMSE | MSE |
|---|---|---|---|---|
| $\widehat{m}(x_1)$ | 0.6082211 | 0.05766383 | 0.08078553 | 0.006526302 |
| $\widehat{m}(x_2)$ | 0.4272434 | 0.1169859 | 0.1631432 | 0.02661572 |
| $\widehat{m}(x_3)$ | 0.6979411 | 0.07678413 | 0.09704057 | 0.009416873 |

Source/ From the (R.4.1.2) Package

Explanation of applied analysis

The explanatory variables that achieved the lowest value for the MSE criterion, respectively (X1= 0.006526302 , X3=0.009416873, X2 =0.02661572) As for the coefficient of determination, the explanatory variables that reached the highest coefficient of determination That is, the explanatory variable explains a percentage of the changes that occur in the dependent variable and the rest is left to other factors , are ( X3 ,X1 ,X2) That is, electrical contact and children's tampering are the main causes of fires .



الدالة $\widehat{m}(x_1)$



الدالة $\widehat{m}(x_2)$

الدالة $\hat{m}(x_3)$

The figure shows the behavior of the real data and how close the estimated values are to the real values.

## CONCLUSIONS

1-When implementing simulation experiments using three sample sizes (n = 30, 60, 100) and with a frequency of 500 for each experiment and depending on the comparison criteria at a level of variance ($\sigma^2$= 0.5), it was found that the best estimator of the first nonparametric model is that the estimator (ALLS) is the best estimator, then It is followed by the estimator (KNN).

2-But when the level of variance is ($\sigma^2$=1) at the sample size (n=30), the best estimator for the first model is (ALLS), followed by (KNN) estimator, as the values of the criteria (RMSE), (MSE) and (MAE) are less. With increasing sample sizes and for all estimators used, and increasing the values of (RMSE), (MSE) and (MAE) for all estimators with increasing values of residual variance.

3-Finally, we can say that the best estimation adaptive method for the three criteria and by increasing the sample sizes at three different levels of variance was the ALLS method, which represents the smoothing of the adaptive local polynomial regression.

4- As for the applied analysis, it was found that The explanatory variables that achieved the lowest value for the MSE criterion, respectively (X1, X3, X2) , As for the coefficient of determination, the explanatory variables that reached the highest coefficient of determination are ( X3 ,X1 ,X2) That is, electrical contact and children's tampering are the main causes of fires .

## ACKNOWLEDGMENTS

## REFERENCES

1-Ali, Noor Abdul-Karim (2021) "Comparing some of the traditional and immunized nonparametriccapabilities of the nonparametric regression model with application", Master's thesis, College of Administration and Economics, University of Basra

2-Aydin , D. & Memmedli2 , M. & Omay2, R. (2013) " Smoothing Parameter Selection for Nonparametric Regression Using Smoothing Spline"  EUROPEAN JOURNAL OF PURE AND APPLIED MATHEMATICS Vol. 6, No. 2, 222-238 ISSN 1307-5543

3-Breunig, C. & Chen, X. (2022) " simple adaptive estimation of quadratic functionals in nonparametric iv models" mathematics statistics theory , version, v2

4-Fan, J. & Gijbels, L. (2017 ) " Local Polynomial Modelling and Its Applications ", MONOGRAPHS ON STATISTICS AND APPLIED PROHABILITY , First edition 1996, Taylor & Francis Group journal

5-Chagnya, G. & Roche, A. (2016) " Adaptive estimation in the functional nonparametric regression model " Journal of Multivariate Analysis ,Volume 146, April

6-Hill ,N. J.  ,  Padmanabhan, A. R.  , Puri, M. L. (1988) " Adaptive Nonparametric Procedures and Applications " Journal of the Royal Statistical Society. Series C (Applied Statistics), Vol. 37, No. 2,1988

7-Lacour, C.& Massart, P. (2016)" Minimal penalty for Goldenshluger–Lepski method" Stochastic Processes and their ApplicationsVolume 126, Issue 12, December

8-Lederer, J.& YU , L.& GAYNANOVA , I. (2019) Oracle Inequalities for High-dimensional Prediction, Vol.25 • No. 2 Bernoulli Society for Mathematical Statistics and Probability.

9-Matta, Nour Sabah and Al-Safawi, Safa Younes (2011) "Estimation of nonparametric regression functions using some improvement", Iraqi Journal of Statistical Sciences, Vol. 2011, No. 20

10-Muhammad, Abdul-Hussain Muhammad (2011) "Using the estimator ((Nadaraya-Watson)) kernel in estimating the nonparametric regression function, Al-Qadisiyah Journal for Administrative and Economic Sciences, Volume 13, Issue 1

11-O'Gorman, T. W. (2004) Applied Adaptive Statistical Methods Tests of Significance and Confidence Intervals , by the American Statistical Association and the Society for Industrial and Applied Mathematics.

12-page, S.& grunewalder , S. ( 2021 ) The Goldenshluger–Lepski Method for Constrained Least-Squares Estimators over RKHSs , Journals Project Euclid , Bernoulli 27(4): 2241-2266 , November Volume 27 ,  Issue 4

13-Serdyukova , N. (2012) Adaptive estimation in regression and complexity of approximation of random fields , Statistics Theory (math.ST); Probability (math.PR)

14-Tali, Nada Hussein and Taher, Ahmed Shaker Mohammed Taher (2022) "Using kernel regression in estimating the coefficients of the regression model with random error limits that are self-correlated with practical application", Journal of Administration and Economics - Al-Mustansiriya University, issue 132, pages 234-247

15-Zhang , Z. (2014) Adaptive Robust Regression Approaches in data analysis and their Applications , PhD thesis , University of Cincinnati ,ProQuest Publishing,