

1 This is a draft version of work in progress, content will be revisited in subsequent versions.

2 Robust Local Polynomial Regression with Similarity Kernels

3 Yaniv Shulman
4 *yaniv@shulman.info*

5 Abstract

Local polynomial regression is a powerful and flexible statistical technique that has gained increasing popularity in recent years due to its ability to model complex relationships between variables. Local polynomial regression generalizes the polynomial regression and moving average methods by fitting a low-degree polynomial to a nearest neighbors subset of the data at the location. The polynomial is fitted using weighted ordinary least squares, giving more weight to nearby points and less weight to points further away. Local polynomial regression is however susceptible to outliers and high leverage points which may cause an adverse impact on the estimation accuracy. The main contribution of this paper is to revisit the kernel that is used to produce local regression weights. The simple yet effective idea is to generalize the kernel such that both the predictor and the response are used to calculate weights. Within this framework, two positive definite kernels are proposed that assign robust weights to mitigate the adverse effect of outliers in the local neighborhood by estimating and utilizing the density at the local locations. The method is implemented in the Python programming language and is made publicly available at <https://github.com/yaniv-shulman/rsklpr>. Experimental results on synthetic benchmarks across a range of settings demonstrate that the proposed method achieves competitive results and generally improves on the standard local polynomial regression method.

1. Introduction

Local polynomial regression (LPR) is a powerful and flexible statistical technique that has gained increasing popularity in recent years due to its ability to model complex relationships between variables. Local polynomial regression generalizes the polynomial regression and moving average methods by fitting a low-degree polynomial to a nearest neighbors subset of the data at the location. The polynomial is fitted using weighted ordinary least squares, giving more weight to nearby points and less weight to points further away. The value of the regression function for the point is then obtained by evaluating the fitted local polynomial using the predictor variable value for that data point. LPR has good accuracy near the boundary and performs better than all other linear smoothers in a minimax sense [2]. The biggest advantage of this class of methods is not requiring a prior specification of a function i.e. a parametrized model. Instead only a small number of hyperparameters need to be specified such as the type of kernel, a smoothing parameter and the degree of the local polynomial. The method is therefore suitable for modeling complex processes such as non-linear relationships, or complex dependencies for which no theoretical models exist. These two advantages, combined with the simplicity of the method, makes it one of the most attractive of the modern regression methods for applications that fit the general framework of least squares regression but have a complex deterministic structure.

Local polynomial regression incorporates the notion of proximity in two ways. The first is that a smooth function can be reasonably approximated in a local neighborhood by a simple function such as a linear or low order polynomial. The second is the assumption that nearby points carry more importance in the calculation of a simple local approximation or alternatively that closer points are more likely to interact in simpler ways than far away points. This is achieved by a kernel which produces values that diminish as the distance between the explanatory variables increase to model stronger relationship between closer points.

Methods in the LPR family include the Nadaraya-Watson estimator [10, 18] and the estimator proposed by Gasser and Müller [7] which both perform kernel based local constant fit. These were improved on in terms of asymptotic bias by the proposal of the local linear and more general local polynomial estimators [16, 3, 9, 4, 5]. For a review of LPR methods the interested reader is referred to [2].

LPR is however susceptible to outliers, high leverage points and functions with discontinuities in their derivative which often cause an adverse impact on the regression due to its use of least squares based optimization [17]. The use of unbounded loss functions may result in anomalous observations severely affecting the local estimate. Substantial work has been done to develop algorithms to apply LPR to difficult data. To alleviate the issue [15] employs variable bandwidth to exclude observations for which residuals from the resulting estimator are large. In [3] an iterated weighted fitting procedure is proposed that assigns in each consecutive iteration smaller weights to points that are farther then the fitted values at the previous iteration. The process repeats for a number of iterations and the final values are considered the robust parameters and fitted values. An alternative common approach is to replace the squared prediction loss by one that is more robust to the presence of large residuals by increasing more slowly or a loss that has an upper bound such as the Tukey or Huber loss. These methods however require specifying a threshold parameter for the loss to indicate atypical observations or standardizing the errors using robust estimators of scale [8]. For a recent review of robust LPR and other nonparametric methods see [17, 11]

The main contribution of this paper is to revisit the kernel used to produce regression weights. The simple yet effective idea is to generalize the kernel such that both the predictor and the re-

sponse are used to calculate weights. Within this framework, two positive definite kernels are proposed that assign robust weights to mitigate the adverse effect of outliers in the local neighborhood by estimating the density of the response at the local locations. Note the proposed framework does not preclude the use of robust loss functions, robust bandwidth selectors and standardization techniques. In addition the method is implemented in the Python programming language and is made publicly available. Experimental results on synthetic benchmarks demonstrate that the proposed method achieves competitive results and generally performs better than LOESS/LOWESS using only a single training iteration.

The remainder of the paper is organized as follows: In section 2, a brief overview of the mathematical formulation of local polynomial regression is provided. In section 3, a framework for robust weights as well as specific robust positive definite kernels are proposed. In section 5, implementation notes and experimental results are provided. Finally, in section 6, the paper concludes with directions for future research.

2. Local polynomial regression

This section provides a brief overview of local polynomial regression and establishes the notation subsequently used. Let (X, Y) be a random pair and $\mathcal{D}_T = \{(X_i, Y_i)\}_{i=1}^T \subseteq \mathcal{D}$ be a training set comprising a sample of T data pairs. Suppose that $(X, Y) \sim f_{XY}$ a continuous density and $X \sim f_X$ the marginal distribution of X . Let $Y \in \mathbb{R}$ be a continuous response and assume a model of the form $Y_i = m(X_i) + \epsilon_i$, $i = 1, \dots, T$ where $m(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ is an unknown function and ϵ_i are independently distributed error terms having zero mean representing random variability not included in X_i such that $\mathbb{E}[Y | X = x] = m(x)$. There are no global assumptions about the function $m(\cdot)$ other than that it is smooth and that locally it can be well approximated by a low degree polynomial as per Taylor's theorem. Local polynomial regression is a class of nonparametric regression methods that estimate the unknown regression function $m(\cdot)$ by combining the classical least squares method with the versatility of non-linear regression. The local p -th order Taylor expansion for $x \in \mathbb{R}$ near a point X_i yields:

$$m(X_i) \approx \sum_{j=0}^p \frac{m^{(j)}(x)}{j!} (x - X_i)^j := \sum_{j=0}^p \gamma_j(x) (x - X_i)^j \quad (1)$$

To find an estimate $\hat{m}(x)$ of $m(x)$ the low-degree polynomial (1) is fitted to the N nearest neighbors using weighted least squares such to minimize the empirical loss $\mathcal{L}_{lpr}(\cdot; \mathcal{D}_N, h)$:

$$\mathcal{L}_{lpr}(x; \mathcal{D}_N, h) := \sum_{i=1}^N \left(Y_i - \sum_{j=0}^p \gamma_j(x) (x - X_i)^j \right)^2 K_h(x - X_i) \quad (2)$$

$$\hat{\gamma}(x) := \min_{\gamma(x)} \mathcal{L}_{lpr}(x; \mathcal{D}_N, h) \quad (3)$$

Where $\gamma, \hat{\gamma} \in \mathbb{R}^{p+1}$; $K_h(\cdot)$ is a scaled kernel, $h \in \mathbb{R}_{>0}$ is the bandwidth parameter and $\mathcal{D}_N \subseteq \mathcal{D}_T$ is the subset of N nearest neighbors of x in the training set where the distance is measured on the predictors only. Having computed $\hat{\gamma}(x)$ the estimate of $\hat{m}(x)$ is taken as $\hat{\gamma}(x)_1$. Note the

term kernel carries here the meaning typically used in the context of nonparametric regression i.e. a non-negative real-valued weighting function that is typically symmetric, unimodal at zero, integrable with a unit integral and whose value is non-increasing for the increasing distance between the X_i and x . Higher degree polynomials and smaller N generally increase the variance and decrease the bias of the estimator and vice versa [2]. For derivation of the local constant and local linear estimators for the multidimensional case see [6].

3. Robust weights with similarity kernels

The main idea presented is to generalize the kernel function used in equation (2) to produce robust weights. This is achieved by using a similarity kernel function defined on the data domain $K_{\mathcal{D}} : \mathcal{D}^2 \rightarrow \mathbb{R}_+$ that enables weighting each point and incorporating information on the data in the local neighborhood in relation to the local regression target.

$$\mathcal{L}_{rsk}(x, y; \mathcal{D}_N, H) := \sum_{i=1}^N \left(Y_i - \sum_{j=0}^p \beta_j(x, y)(x - X_i)^j \right)^2 K_{\mathcal{D}}((x, y), (X_i, Y_i); H) \quad (4)$$

$$\hat{\beta}(x, y; \mathcal{D}_N, H) := \min_{\beta(x, y)} \mathcal{L}_{rsk}(x, y; \mathcal{D}_N, H) \quad (5)$$

Where H is the set of bandwidth parameters. There are many possible choices for such a similarity kernel to be defined within this general framework. However, used as a local weighting function, such a kernel should have the following attributes:

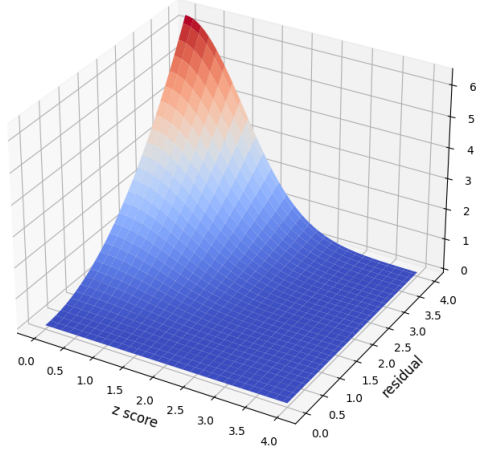
1. Non-negative, $K_{\mathcal{D}}((x, y), (x', y')) \geq 0$.
2. Symmetry in the inputs, $K_{\mathcal{D}}((x, y), (x', y')) = K_{\mathcal{D}}((x', y'), (x, y))$.
3. Tending toward decreasing as the distance in the predictors increases. That is, given a similarity function on the response $s(\cdot, \cdot) : \mathbb{R}^2 \rightarrow \mathbb{R}_+$, if $s(y, y')$ indicates high similarity the weight should decrease as the distance between the predictors grows, $s(y, y') > \alpha \implies K_{\mathcal{D}}((x, y), (x + u, y')) \geq K_{\mathcal{D}}((x, y), (x + v, y')) \quad \forall \|u\| \leq \|v\|$ and some $\alpha \in \mathbb{R}_+$.

In this work two such useful positive definite kernels are proposed. Similarly to the usual kernels used in (2), these tend to diminish as the distance between the explanatory variables increases to model stronger relationship between closer points. In addition, the weights produced by the kernels also model the "importance" of the pair (x, y) . This is useful for example to down-weight outliers to mitigate their adverse effect on the ordinary least square based regression. Formally let $K_{\mathcal{D}}$ be defined as:

$$K_{\mathcal{D}}((x, y), (x', y'); H_1, H_2) = K_1(x, x'; H_1) K_2((x, y), (x', y'); H_2) \quad (6)$$

Where $K_1 : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$ and $K_2 : \mathcal{D}^2 \rightarrow \mathbb{R}_+$ are positive definite kernels and H_1, H_2 are the sets of bandwidth parameters. The purpose of K_1 is to account for the distance between a neighbor to the local regression target and therefore may be chosen as any of the kernel functions that are typically used in equation (2). The role of K_2 is described now in more detail as this is the main idea proposed in this work. Using K_2 , the method performs robust regression by detecting local outliers in an unsupervised manner and assigns them with lower weights. There are many methods that could be employed to estimate the extent to which a data point is a local outlier however in this work it is estimated in one of the following two ways.

Figure 1: Loss function, assuming a standard quadratic function of the residual, a standard normal density for K_2 and excluding the K_1 distance kernel scaling.



117 *Conditional density*

The first proposed method for K_2 is proportional to the estimated localized conditional marginal distribution of the response variable at the location:

$$K_2((x, y), (x', y'); H_2) = \hat{f}(y | x; H_2) \hat{f}(y' | x'; H_2) \quad (7)$$

The nonparametric conditional density estimation is performed using the Parzen–Rosenblatt window (kernel density estimator):

$$\hat{f}(y | x; H_2) = \hat{f}(x, y; H_2) / \hat{f}(x; H_2) \quad (8)$$

$$= \hat{f}(v; \mathbf{H}_v) / \hat{f}(x; \mathbf{H}_x) \quad (9)$$

$$= \frac{|\mathbf{H}_x|^{1/2} \sum_{i=1}^N K_v(\mathbf{H}_v^{-1/2}(v - V_i))}{|\mathbf{H}_v|^{1/2} \sum_{i=1}^N K_x(\mathbf{H}_x^{-1/2}(x - X_i))} \quad (10)$$

118 Where $v = [x, y] \in \mathbb{R}^{d+1}$ is the concatenated vector of the predictors and the response; and $\mathbf{H}_v, \mathbf{H}_x$
119 are bandwidth matrices.

120 *Joint density*

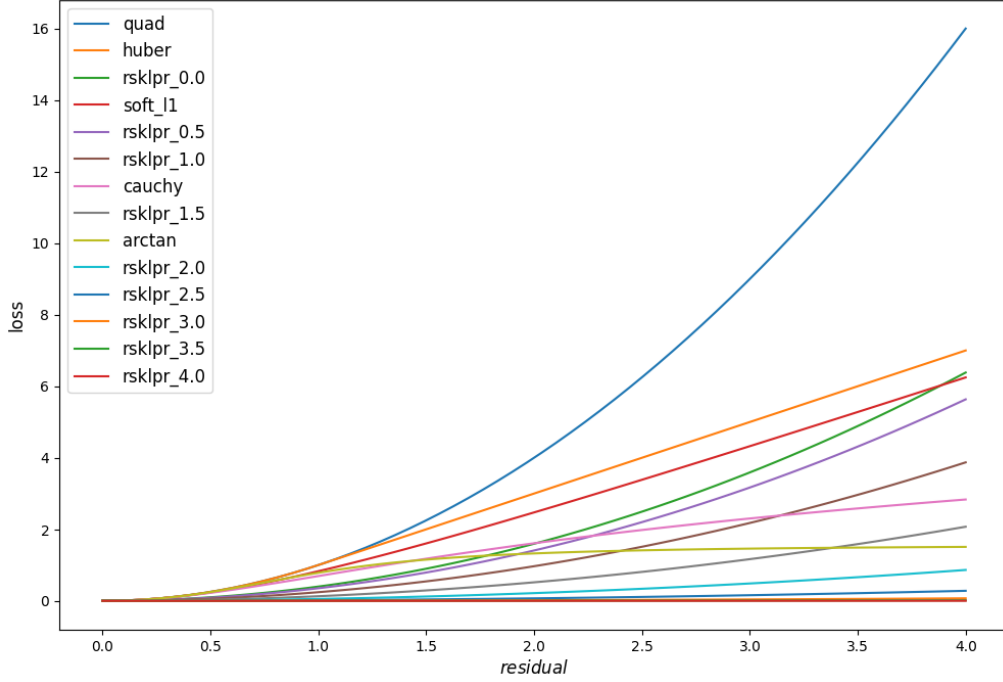
The second proposed kernel is proportional to the joint distribution of the random pair, this could be useful for example to also down-weight high leverage points:

$$K_2((x, y), (x', y'); H_2) = \hat{f}(x, y; H_2) \hat{f}(x', y'; H_2) \quad (11)$$

121 Where the joint density can be estimated using the same aforementioned approach.

122

Figure 2: The plot illustrates the proposed loss function, a number of common robust losses and the standard quadratic residual loss for comparison. It is assumed that that K_2 is equivalent to the standard normal density and the K_1 distance kernel scaling is excluded. The numbers appended to "rsklpr" indicate how many standard deviations away from the mean the density is calculated. It is evident that the loss is heavily attenuated in regions of low density.



1

123 Regardless of the choice of kernel, the hyperparameters of this model are similar in essence
 124 to the standard local polynomial regression and comprise the span of included points, the kernels
 125 and their associated bandwidths. Note that this estimator can be replaced with other robust
 126 density estimators and better results are anticipated by doing so however exploring this option is
 127 left for future work.

128 4. Properties

129 This section discusses some properties of the estimator. Note the notation in this section is
 130 simplified by excluding explicit mentions of D_N and H , however the analysis is conditional on
 131 the nearest neighbors in the sample, D_N .

132 4.1. Invariance to y at the regression location and simplification of the objective

133 The objective (5) is invariant to the value of y at the location (x, y) for the proposed similarity
 134 kernels.

135 *Proof:* The optimization is invariant to the scale of the objective function. Therefore:

$$\hat{\beta}(x, y) := \min_{\beta(x, y)} \sum_{i=1}^N \left(Y_i - \sum_{j=0}^p \beta_j(x, y)(x - X_i)^j \right)^2 K_{H_1}(x - X_i) \hat{f}(x, y) \hat{f}(X_i, Y_i) \quad (12)$$

$$= \sum_{i=1}^N \left(Y_i - \sum_{j=0}^p \beta_j(x, y)(x - X_i)^j \right)^2 K_{H_1}(x - X_i) \hat{f}(X_i, Y_i) \quad (13)$$

136 The equality holds because $\hat{f}(x, y)$ is a constant scalar that uniformly scales the weights.
 137 Since the objective is now independent of y , it follows that:

$$\hat{\beta}(x, y) := \min_{\beta(x)} \sum_{i=1}^N \left(Y_i - \sum_{j=0}^p \beta_j(x)(x - X_i)^j \right)^2 K_{H_1}(x - X_i) \hat{f}(X_i, Y_i) \quad (14)$$

$$:= \hat{\beta}(x) \quad \forall y \quad (15)$$

138 This simplification enables more efficient calculations of the estimator because the depen-
 139 dence on y is removed from the objective function. Note that $\hat{f}(X_i, Y_i)$ can also be replaced with
 140 $\hat{f}(Y_i | X_i)$ with similar results.

141 4.2. Weighted arithmetic mean of the standard LPR

142 The proposed estimator is equivalent to the weighted arithmetic mean of the terms in the
 143 standard LPR loss (2), with weights $w_i = \hat{f}(X_i, Y_i)$.

144 *Proof:* Since the optimization is invariant to scaling, we have:

$$\hat{\beta}(x) := \min_{\beta(x)} \sum_{i=1}^N \left(Y_i - \sum_{j=0}^p \beta_j(x)(x - X_i)^j \right)^2 K_{H_1}(x - X_i) \hat{f}(X_i, Y_i) \quad (16)$$

$$= \min_{\beta(x)} \left(\sum_{i=1}^N \hat{f}(X_i, Y_i) \right)^{-1} \sum_{i=1}^N \left(Y_i - \sum_{j=0}^p \beta_j(x)(x - X_i)^j \right)^2 K_{H_1}(x - X_i) \hat{f}(X_i, Y_i) \quad (17)$$

$$= \min_{\beta(x)} \left(\sum_{i=1}^N w_i \right)^{-1} \sum_{i=1}^N \left(Y_i - \sum_{j=0}^p \beta_j(x)(x - X_i)^j \right)^2 K_{H_1}(x - X_i) w_i \quad (18)$$

145 The normalization by $\sum_{i=1}^N w_i$ shows the equivalence to the weighted arithmetic mean, with
 146 the weights $w_i = \hat{f}(X_i, Y_i)$.

147 4.3. Asymptotic degeneration of the estimator to the standard LPR

148 Asymptotically, the proposed estimator degenerates to the standard LPR when the weights
 149 w_i are uncorrelated with the standard LPR terms. Formally, as $N \rightarrow \infty$, $\hat{\beta}(x) \rightarrow \hat{\gamma}(x)$, where
 150 $\hat{\gamma}(x)$ is the standard LPR estimator, and the condition $\left(Y - \sum_{j=0}^p \beta_j(x)(x - X)^j \right)^2 K_{H_1}(x - X)$ and

¹⁵¹ $\hat{f}(X, Y)$ are uncorrelated. We assume that (X_i, Y_i) are independent and identically distributed
¹⁵² (i.i.d.) random variables and that $\hat{f}(X, Y) > 0$ almost everywhere.

¹⁵³ *Proof:* Define

$$g(X, Y) := \left(Y - \sum_{j=0}^p \beta_j(x)(x - X)^j \right)^2 K_{H_1}(x - X),$$

it follows that:

$$\hat{\beta}(x) := \min_{\beta(x)} \left(\sum_{i=1}^N \hat{f}(X_i, Y_i) \right)^{-1} \sum_{i=1}^N \left(Y_i - \sum_{j=0}^p \beta_j(x)(x - X_i)^j \right)^2 K_{H_1}(x - X_i) \hat{f}(X_i, Y_i) \quad (19)$$

$$= \min_{\beta(x)} \left(\sum_{i=1}^N \hat{f}(X_i, Y_i) \right)^{-1} \sum_{i=1}^N g(X_i, Y_i) \hat{f}(X_i, Y_i) \quad (20)$$

$$= \min_{\beta(x)} \left(\frac{1}{N} \sum_{i=1}^N \hat{f}(X_i, Y_i) \right)^{-1} \frac{1}{N} \sum_{i=1}^N g(X_i, Y_i) \hat{f}(X_i, Y_i) \quad (21)$$

As $N \rightarrow \infty$, by the law of large numbers and the continuous mapping theorem, we obtain:

$$\lim_{N \rightarrow \infty} \left(\frac{1}{N} \sum_{i=1}^N \hat{f}(X_i, Y_i) \right)^{-1} \xrightarrow{a.s.} \mathbb{E} \left[\frac{1}{\hat{f}(X, Y)} \right] \quad (22)$$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N g(X_i, Y_i) \hat{f}(X_i, Y_i) \xrightarrow{a.s.} \mathbb{E} [g(X, Y) \hat{f}(X, Y)] \quad (23)$$

Assuming $\mathbb{E}[\hat{f}(X, Y)] \neq 0$, it follows that:

$$\hat{\beta}(x) := \min_{\beta(x)} \frac{\mathbb{E} [g(X, Y) \hat{f}(X, Y)]}{\mathbb{E} [\hat{f}(X, Y)]} \quad (24)$$

If $g(X, Y)$ and $\hat{f}(X, Y)$ are uncorrelated, then:

$$\hat{\beta}(x) := \min_{\beta(x)} \frac{\mathbb{E} [g(X, Y) \hat{f}(X, Y)]}{\mathbb{E} [\hat{f}(X, Y)]} = \min_{\beta(x)} \mathbb{E} [g(X, Y)] \quad (25)$$

Thus, as $N \rightarrow \infty$:

$$\hat{\beta}(x) = \min_{\beta(x)} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N g(X_i, Y_i) \quad (26)$$

$$= \min_{\beta(x)} \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \left(Y_i - \sum_{j=0}^p \beta_j(x)(x - X_i)^j \right)^2 K_{H_1}(x - X_i) \quad (27)$$

¹⁵⁴ By removing the $\frac{1}{N}$ term from the summation in the limiting case, we recover the standard
¹⁵⁵ LPR estimator. This is because the optimization is invariant to scale, and removing $\frac{1}{N}$ simply

156 leads to a minimization of the sum of squared residuals, which is the essence of the standard
157 LPR. Thus, the proposed estimator degenerates to the standard LPR in this finite-sample setting
158 as $N \rightarrow \infty$, provided the sum does not diverge.

159 The goal of the proposed method is to prevent the degeneration to standard LPR by ensuring
160 that the weights $f(X, Y)$ are correlated with the residuals. If $f(X, Y)$ becomes uncorrelated with
161 $g(X, Y)$, the weighting scheme would lose its intended robustness, treating all points equally,
162 and the estimator would revert to the standard LPR, losing the benefits of outlier mitigation.
163 The method is designed to avoid this by using similarity kernels that ensure the weight func-
164 tion $f(X, Y)$ reflects the local structure of the data, with the expectation that large residuals are
165 indicative of outliers or lower density regions.

166 5. Experiments and implementation notes

167 The proposed method was implemented in Python. The distances between pairs are normal-
168 ized to the range $[0, 1]$ in each neighborhood as in [3]. The simple Laplacian kernel $e^{-\|x-x'\|}$ is
169 used for $K_1(x, x'; H_1)$ since it gave better and more efficient empirical results in tests than the
170 tricube kernel recommended in [3]. A factorized multidimensional KDE using scaled Gaussian
171 kernels is used for estimating the density. Five methods for estimating the bandwidth were used:
172 Scott’s rule [12], Normal Reference, global LSCV, local LSCV and local MLCV. Additionally
173 the bandwidth for the predictor’s kernel also uses in some of the experiments a simple function of
174 the window size estimated empirically. Certain computations are omitted for efficiency since the
175 local regression in equation (5) is invariant to the scale of the weights. This includes all scaling
176 constants fixed in a given neighborhood concerning a specific local regression target including
177 the computation of $\hat{f}(y|x)$ and $\hat{f}(x, y)$ in equations (7) and (11) respectively. Experiments were
178 performed on a number of synthetic benchmarks to evaluate the effectiveness of the method’s
179 linear and quadratic variants in comparison to other methods of the local polynomial regression
180 family including LOWESS and iterative robust LOWESS [13], local linear and local constant
181 kernel regression [13], local quadratic regression [14] and radial basis function network [14].

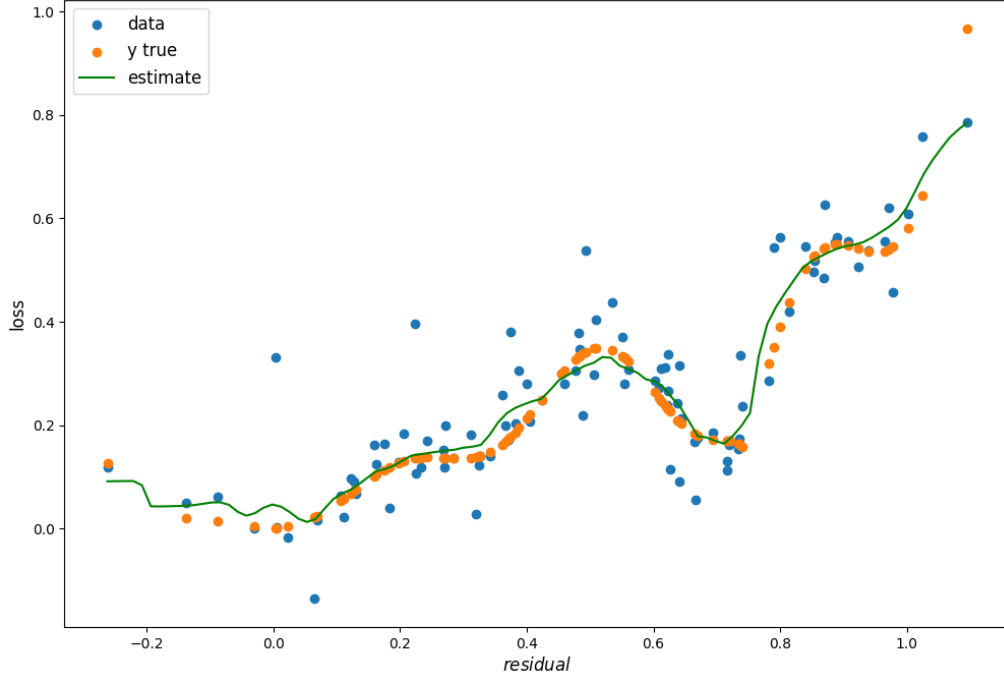
182
183 The experiments comprise a number settings including various non-linear synthetic curves
184 and planes with added noise. These have variants of dense and sparse data, homoscedastic and
185 heteroscedastic noise characteristics and various neighborhood sizes. The results indicate there
186 is no best universal method and that different methods work better in a given setting. However,
187 the proposed method offers competitive performance across the board and gives the best results
188 in a large number of settings and in particular in heteroscedastic settings. It is further shown the
189 proposed method generally improves on it’s direct counterparts LOESS/LOWESS and quadratic
190 LPR using a single iteration. In addition it appears far less sensitive to the choice of neighborhood
191 size and has substantial reduced variance in this respect. This last quality makes it an attractive
192 choice since it is more likely for the analyst to select an “appropriate” hyperparameter value
193 when there is no ground truth available.

194 The experimental results are available at [https://nbviewer.org/github/yaniv-shulman/](https://nbviewer.org/github/yaniv-shulman/rsklpr/tree/main/src/experiments/)
195 [rsklpr/tree/main/src/experiments/](https://nbviewer.org/github/yaniv-shulman/rsklpr/tree/main/src/experiments/) as interactive Jupyter notebooks [1].

196 6. Future work and research directions

197 This work proposes a new robust variant of LPR and as such there are many aspects that need
198 to be explored and are left for future work. As the proposed method generalize the standard LPR

Figure 3: Regression example of synthetically generated 1D data with heteroscedastic noise. Additional experimental results and demonstrations including multivariate settings and bootstrap based confidence intervals are available at <https://nbviewer.org/github/yaniv-shulman/rsklpr/tree/main/src/experiments/> as interactive Jupyter notebooks [1]



it does not preclude replacing some of the standard LPR components in equation (5) with other and robust alternatives including robust methods for bandwidth selection and robust alternatives to the standard quadratic residual function. Another avenue for developing this framework is investigating additional kernels K_D and exploring the impact of robust density estimators on performance.

References

- [1] Project jupyter is a non-profit, open-source project, born out of the ipython project in 2014 as it evolved to support interactive data science and scientific computing across all programming languages. <https://jupyter.org/>.
- [2] M. Avery. Literature review for local polynomial regression. 2010.
- [3] W. S. Cleveland. Robust locally weighted regression and smoothing scatterplots. *Journal of the American Statistical Association*, 74(368):829–836, 1979.
- [4] W. S. Cleveland and S. J. Devlin. Locally weighted regression: An approach to regression analysis by local fitting. *Journal of the American Statistical Association*, 83(403):596–610, 1988.
- [5] J. Fan. Local linear regression smoothers and their minimax efficiencies. *The Annals of Statistics*, 21, 03 1993.
- [6] E. García-Portugués. *Notes for Nonparametric Statistics*. 2023. Version 6.9.0. ISBN 978-84-09-29537-1.
- [7] T. Gasser and H.-G. Müller. Estimating regression functions and their derivatives by the kernel method. *Scandinavian Journal of Statistics*, 11:171–185, 1984.
- [8] R. A. Maronna, D. Martin, V. J. Yohai, and Hardcover. Robust statistics: Theory and methods. 2006.
- [9] H.-G. Muller. Weighted local regression and kernel methods for nonparametric curve fitting. *Journal of the American Statistical Association*, 82(397):231–238, 1987.

- 219 [10] E. Nadaraya. On estimating regression. *Theory of Probability and Its Applications*, 9:141–142, 1964.
- 220 [11] M. Salibian-Barrera. Robust nonparametric regression: Review and practical considerations. *Econometrics and*
- 221 *Statistics*, 2023.
- 222 [12] D. Scott. *Multivariate Density Estimation: Theory, Practice, and Visualization*. Wiley Series in Probability and
- 223 *Statistics*. Wiley, 2015.
- 224 [13] S. Seabold and J. Perktold. statsmodels: Econometric and statistical modeling with python. In *9th Python in*
- 225 *Science Conference*, 2010.
- 226 [14] sigvaldm. Localreg is a collection of kernel-based statistical methods. [https://github.com/sigvaldm/](https://github.com/sigvaldm/localreg)
- 227 [localreg](https://github.com/sigvaldm/localreg).
- 228 [15] V. G. Spokoiny. Estimation of a function with discontinuities via local polynomial fit with an adaptive window
- 229 *choice*. *The Annals of Statistics*, 26(4):1356 – 1378, 1998.
- 230 [16] C. J. Stone. Consistent nonparametric regression. *Annals of Statistics*, 5:595–620, 1977.
- 231 [17] P. Čížek and S. Sadıkoğlu. Robust nonparametric regression: A review. *WIREs Comput. Stat.*, 12(3), apr 2020.
- 232 [18] G. S. Watson. Smooth regression analysis. 1964.