

Iterative Reweighted Quantile Regression Using Augmented Lagrangian Optimization for Baseline Correction

Quanjie Han, Silong Peng*, Qiong Xie, Yifan Wu
Institute of Automation, Chinese Academy of Sciences
University of Chinese Academy of Sciences
Beijing 100190, PR China
hanquanjie2015@ia.ac.cn, silong.peng@ia.ac.cn,
qiong.xie@ia.ac.cn, yfwu5216@ia.ac.cn

Genwei Zhang
Beijing Institute of Pharmaceutical Chemistry
Beijing 102205, PR China
zhgw1984@163.com

Abstract—Based on baseline is a smooth curve and under the collected spectrum, a robust penalized quantile regression with B-spline basis has been proposed to baseline estimation. Then an iterative reweighted method has been adopted for quantile regression optimization. Instead of man tuning the hyperparameter in penalized quantile regression, augmented Lagrangian method is applied to hyperparameter optimization. Experiments on simulated and real data sets show that our method is more effective in baseline correction than other baseline estimation methods in simulated data set. For real data set, the calibration results after the baseline correction step are better than other preprocessing and baseline correction methods.

Keywords—quantile regression; p-splines; iterative reweighted least squares; augmented Lagrangian

I. INTRODUCTION

Since Fourier transform spectrometer is rapid and nondestructive, Fourier transform infrared spectroscopy has been widely used in Chemometrics, food, wine and other related fields for sample components analysis [1]. Generally speaking, the obtained Fourier transform infrared absorption spectroscopy consists of the true sample spectrum, baseline and noise. Baseline together with noise will significantly deteriorate the performance of chemometric calibration algorithms, so baseline correction and spectrum denoising is an important preprocessing step for spectrum quantitative analysis.

Baseline estimation can be dated back to late 1970s [2]. Up to now, there have several assumptions been imposed on baseline: Firstly, from the frequency prospective, baseline is in low frequency part while noise generally lives in the high frequency part, low-pass filter has been constructed to correct the baseline [3, 4]. Secondly, baseline is a smooth curve which underlies the collected spectrum: it was fitted by polynomials [5] and Bernstein polynomials were proposed for extraction of baseline of NMR signals [6]. Last but not the least, baseline points and spectrum peak points belong to different clusters, which can be separated. In order to separate the baseline points and peak points, Rooi proposed a mixture model for baseline estimation. The baseline points

were characterized by a Gaussian distribution while the peak points were subject to a uniform distribution. Using EM algorithm, after the baseline points and peak points having been separated, a penalized B-splines basis was used to fitting the baseline [7]. The Gaussian mixture model was also used for DNA sequence baseline correction in [8]. Since baseline underlies the obtained spectrum, an asymmetrically weighted least squares (asLS) with roughness penalty was proposed for baseline estimation [9]. The weights for the baseline points below the spectrum were set manually by a constant which usually will overestimate the peak and there were two parameters need to be optimized, [10] proposed the adaptive iteratively reweighted Penalized Least Squares (airPLS) and a partially balanced weighting scheme was also proposed in [11] for baseline smoothing (arPLS). Besides, based on the spectrum of the sample can be approximated by Voight lineshape, a method simultaneously fitting the pure spectrum and baseline using sparse representation (SSFBCSP) was proposed in [12] and a multiple spectral baseline correction method which combined the information of several spectral was used for Guotai wine baseline correction [13]. Simple Least squares regression corresponding to the conditional mean value regression and the least squares regression is sensitive to outliers and noise. Besides, in order to obtain the regression equation for other quantiles, quantile regression was proposed by Koenker and Bassett [14] and it was first used for baseline correction in [15].

This article proposes a quantile regression with penalized B-splines for baseline correction, where the B-splines are used to represent the baseline. Instead of the primal-dual interior or simplex method for solving quantile regression problem, we propose the iterative reweighted least squares to tackle the quantile regression, which has several advantages: it is easy to implement than the linear programming methods; iterative reweighted least squares usually gives more accurate result. In order to avoid the optimization of the regularization parameter in penalized B-splines, the augmented Lagrangian method is also proposed for hyperparameter optimization. This paper is divided into the

following parts: Section II provides the detailed introduction for our algorithms; Section III displays our experiments setting; The experiments results and discussion are shown in Section IV; The final part is devoted to Conclusion.

II. PROBLEM FORMULATION

A. P-splines and Quantile Regression

Roughness penalty approach to problems in regression has gained much popularity in recent years, especially in functional data analysis (FDA) [16]. Considering the nonparametric regression problem: given n data pairs $(x_i, y_i), i = 1, 2, \dots, n$, find a function g such that

$$y_i = g(x_i) + e_i \quad (1)$$

Where e_i is the error term in i -th sample, which is usually assumed normally distributed. Without constraint on g , then the error term can be zero just by interpolating the points using piecewise linear function. In order to find a compromise between the fidelity of curve fitting and avoiding of rapidly fluctuating curve, a penalty is posed on the curvature of g , then the smooth penalty based regression becomes

$$\sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int g''(x) dx \quad (2)$$

the second derivative can be replaced by other higher order derivatives. For computation convenience, the function g is represented by some basis functions usually. In functional data analysis, the Fourier basis is used for periodic functions, while B-spline basis is adopted for nonperiodic functions. Assume that g can be represented by B-spline basis

$$g = B\alpha = \sum_j B_{i,j}\alpha_j \quad (3)$$

By the recurrence relation and formula for derivatives of B-splines given by de Boor [17], the continuous penalty $\int g''(x) dx$ is equivalent to the smooth of the representation coefficient. Then (2) becomes

$$\|y - B\alpha\|_2^2 + \lambda \|D\alpha\|_2^2 \quad (4)$$

where D is the difference matrix, the difference order is usually set to and three. The B-splines with penalty is called P-splines [18].

From the maximal likelihood point of view, the least squared term is corresponding to the noise is Gaussian and what we obtain is the conditional expectation of y given x , which is sensitive to outliers. To get robust estimator, the absolute deviation has been adopted which corresponding to the conditional median. In order to get the information of other quantiles, Koenker proposed the quantile regression, which can be formulated as an optimization problem:

$$\arg \min_z \sum_i \rho_\tau(y_i - z_i) \quad (5)$$

where $\rho_\tau = u(\tau - 1(u < 0)) = \tau u_+ + (1 - \tau)u_-$, $u_+ = \max\{u, 0\}$ is the positive part of u , while $u_- = \max\{-u, 0\}$ is the negative part of u . The median regression is corresponding to $\tau = 0.5$.

B. Iterative Reweighted Quantile Regression With Augmented Lagrangian Optimization

Since baseline is running below the spectrum, we should impose asymmetrically penalty for estimated points. For points above the original spectrum, a large penalty should be set, while for points under it, a small penalty should imposed. From quantile point of view, the baseline is at the low quantile part of the original spectrum.

Due to the smooth of baseline, it can be represented by B-splines

$$z = B\alpha \quad (6)$$

Then the quantile regression with P-splines for baseline correction can be described as follows:

$$\arg \min_\alpha \sum_i \rho_\tau(y_i - \sum_j B_{i,j}\alpha_j) + \lambda \|D\alpha\|^2 \quad (7)$$

After the representation coefficient α is obtained, then the baseline is estimated by $z = B\alpha$.

Considering that $\rho_\tau = u(\tau - 1(u < 0)) = \tau u_+ + (1 - \tau)u_-$ is not differentiate at zero. By using iterative reweighted least squares and noting that $|u| = \frac{u^2}{|u|}$, we can set the asymmetric weight as

$$w_i = \begin{cases} \frac{\tau}{|y_i - z_i| + \epsilon} & y_i \geq z_i \\ \frac{1 - \tau}{|y_i - z_i| + \epsilon} & y_i < z_i \end{cases} \quad (8)$$

where $\epsilon > 0$ is added by avoiding the divided by zero problem. Then the quantile regression can be optimized by

$$\arg \min_\alpha \sum_i w_i (y_i - \sum_j B_{i,j}\alpha_j)^2 + \lambda \|D\alpha\|_2^2 \quad (9)$$

In reality, the success of baseline estimation depends on choosing the hyper-parameter λ properly. To avoid the tuning of λ , we propose to paraphrase (9) as

$$\arg \min_\alpha \sum_i w_i (y_i - \sum_j B_{i,j}\alpha_j)^2, -\epsilon \leq D\alpha \leq \epsilon \quad (10)$$

where the inequality is applied element-wise. With augmented Lagrangian optimization, (10) becomes

$$\arg \min_\alpha \sum_i w_i (y_i - \sum_j B_{i,j}\alpha_j)^2 + v^T D\alpha + \frac{\rho}{2} \|D\alpha\|_2^2 \quad (11)$$

Let W denote the diagonal matrix with $w = (w_i)$ on its diagonal, (11) can be described as

$$\arg \min_\alpha (y - B\alpha)^T W (y - B\alpha) + v^T D\alpha + \frac{\rho}{2} \|D\alpha\|_2^2 \quad (12)$$

where v is the Lagrangian multipliers and ρ is a penalty parameter. We summarize the iterative reweighted quantile

regression with augmented Lagrangian Optimization method (IRQRAL) as follows:

Algorithm: IRQRAL

Step 1. Input single spectrum y , penalty parameter ρ , quantile τ , order of difference matrix d , B-spline basis matrix B , maximum iteration number $Iter$.

Step 2. Initialize $w = 1_n$, $k = 0$, ρ_{\max} , ϵ , relative error ϵ_1 .

Step 3. Update α , v and ρ :

3.1 $W = \text{diag}(w)$;

3.2 $\alpha^{(k+1)} = (2B^T W B + \rho D^T D)^{-1} (2B^T W y - D^T v)$;

3.3 $v^{(k+1)} = v^{(k)} + \rho D \alpha^{(k+1)}$.

3.4 $\rho = \min(2\rho, \rho_{\max})$.

Step 4 Reweight w with (8).

Step 5. Check stopping criterion:

if $\|\alpha^{(k)} - \alpha^{(k-1)}\| < \epsilon_1$ or $k > Iter$, stop; else $k \leftarrow k+1$ go to **Step 3**.

Step 6. Output baseline $z = B\alpha$, representation coefficient α .

III. EXPERIMENTS

A. Dataset description

To evaluate the performance of the proposed method, one simulated data and one real data set are used for quantitative analysis. The simulated data consists of six Gaussian peaks and a sinusoidal baseline and an exponential baseline are added to the peaks respectively. Besides, a uniform random noise is generated whose amplitude doesn't exceed 0.01 of the maximal height of the peaks. The real data is the corn data set which consists of 80 NIR spectra of corn measured on spectrometers mp5 and mp6 respectively and the spectra were collected in the region of 1100-2498 nm. There are four components have been measured: moisture, oil, protein, starch. In our experiment, the mp5 data set is adopted to compare the calibration result of our method with other preprocessing methods after the baselines being corrected.

B. Model evaluation

For simulated data set, the true baseline is known, we can use root mean squared error (RMSE) to find the optimal parameters. In our experience, the quantile $\tau = 0.01$ and the order of difference matrix $d = 3$ can always give desired results. Since whatever we choose ρ , the IRQRAL algorithm will converge, so we fix $\rho = 1$. The RMSE is computed by

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (z_i - \hat{z}_i)^2}{n}} \quad (13)$$

where z is the true baseline, \hat{z} is the estimated baseline and n is the length of simulated data.

Since the true baselines for real data set are unavailable, the quantitative results of spectra after the baselines having been corrected are used to metric the success of baselines

estimation. We split the data set into training set and test set. Firstly, each response component is sorted, then the second of every four is taken as test set, the others are treated as training set; then leaving one out cross validation is used for calibration. In order to avoid over fitting, a criterion based on testing the significance of incremental changes in PRESS with an F-test [19] is used for the choice of the number of latent variables. In this work, a 95% confidence interval is employed. Finally, the root mean squared error of prediction (RMSEP) is used to evaluate the performance of each method.

IV. RESULTS AND DISCUSSIONS

A. Simulated data sets

With respect to simulated data with sinusoidal baseline and exponential baseline, the estimated baselines by our algorithm are shown in Figure 1. The asLS, airPLS, arPLS, SSFBCSP baseline correction methods and the iterative reweighted quantile regression for baseline estimation without augmented Lagrangian optimization (IRQR) are used to compared with our method, the parameters of each method are optimized by grid search. The RMSE for each baseline correction method are detailed in Table I and Table II respectively.

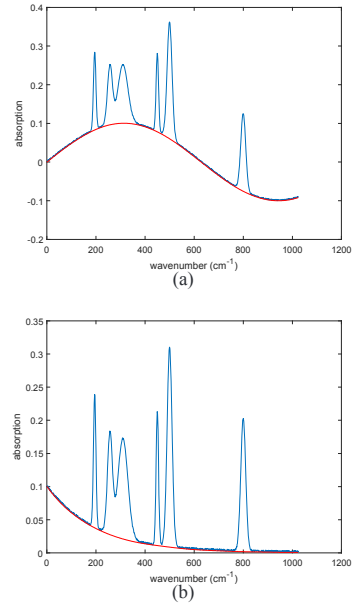


Figure 1. (a) Original spectrum (blue) and estimated baseline (red) for peaks with sinusoidal baseline. (b) Original spectrum (blue) and estimated baseline (red) for peaks with exponential baseline.

We can see that our method is better than the other methods and quantile regression with augmented Lagrangian optimization is generally outperforms the one without it.

Table I
RMSE FOR EACH ESTIMATION METHOD OF SINUSOIDAL BASELINE

Methods	Optimal Parameters	RMSE
asLS	$\lambda = 10^4$ $p = 10^{-6}$	0.0043
airPLS	$\lambda = 10^6$	0.0015
arPLS	$\lambda = 10^5$ $p = 10^{-3}$	0.0019
SSFBCSP	$\lambda_1 = 10^6$ $\lambda_2 = 0.01$	7.83×10^{-4}
IRQR	$\lambda = 10^{11}$	3.4458×10^{-4}
IRQRAL	$\tau = 0.01$	2.6320×10^{-4}

Table II
RMSE FOR EACH ESTIMATION METHOD OF EXPONENTIAL BASELINE

Methods	Optimal Parameters	RMSE
asLS	$\lambda = 10^3$ $p = 10^{-5}$	0.0026
airPLS	$\lambda = 10^6$	0.0011
arPLS	$\lambda = 10^5$ $p = 10^{-3}$	0.0017
SSFBCSP	$\lambda_1 = 10^6$ $\lambda_2 = 0.01$	5.55×10^{-4}
IRQR	$\lambda = 10^{12}$ $\tau = 0.01$	3.2574×10^{-4}
IRQRAL	$\tau = 0.01$	1.9131×10^{-4}

B. Corn data set

The original spectral and the estimated baselines, the baseline corrected spectral by our algorithm are displayed in Figure 2. In Figure 2(a), we can see that corn data set spectral have severe baseline drift, which will adversely influence the calibration and prediction results conducted on it. While seeing from Figure 2(b), our method has successfully corrected the spectral to zero baseline.

To evaluate the performance of our method, spectral preprocessing methods include multiplicative scatter correction (MSC) [20], standard normal variate (SNV) [21], extended inverse scatter correction (EISC) [22], extended multiplicative signal correction (EMSC) [23] and baseline correction methods consist of asymmetrically least squares (asLS), multiple spectral baseline correction (MSBC) and simultaneous spectrum fitting and baseline correction using sparse representation (SSFBCSP) are used to get transformed spectral. The transformed spectra are mean centered before calibration. The RMSEP for moisture, oil, protein, starch are detailed in Table III. For oil component, no methods give desired results. But our algorithm gets better results in other three components and is better than the method without augmented Lagrangian optimization.

C. Discussion

Since our quantile regression problem is just an iterative weighted least squares problem, the Augmented Lagrangian method which can update the dual variables and penalty parameter gradually and the experiments show the desired result. Besides, for the generation of B-splines matrix, the location of knots of B-splines can be set the same as wave number of spectrum, but it may generate many redundant B-spline basis to represent the baseline. Set the location

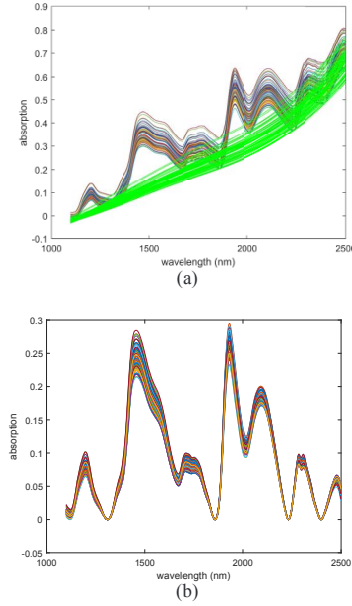


Figure 2. (a) The estimated baselines (green) and original spectra. (b) Baselines corrected spectra.

Table III
RMSEP FOR EACH METHOD

	Moisture	Oil	Protein	Starch
NO	0.1280(6) *	0.0718(5)	0.1778(10)	0.3654(8)
MSC	0.1467(6)	0.0855(4)	0.1812(8)	0.4102(6)
SNV	0.1466(6)	0.0855(4)	0.1813(8)	0.4096(6)
MSBC	0.1203(8)	0.0932(8)	0.1305(15)	0.3359(8)
asLS	0.0957(6)	0.0839(6)	0.1741(11)	0.3292(7)
EISC	0.1730(5)	0.0929(3)	0.1883(4)	0.4256(3)
EMSC	0.1683(5)	0.0948(3)	0.1883(3)	0.4072(3)
SSFBCSP	0.1108(9)	0.0944(7)	0.1303(8)	0.3329(8)
IRQR	0.0804(8)	0.0849(4)	0.1441(9)	0.3261(7)
IRQRAL	0.0763(8)	0.0876(7)	0.1289(9)	0.3201(8)

* The values in parentheses refer the number of latent variables.

of knots optimally, which can improve the computation efficiency.

V. CONCLUSION

The proposed method uses the iterative reweighted quantile regression and augmented Lagrangian optimization for the baseline estimation, which can free us of the parameter optimization used by other baseline estimation methods. And since the quantile regression is more robust than the least squares, this baseline estimation method can not only be used for the Gaussian noise situation, but also for other heterogeneous and dependent data error circumstance.

ACKNOWLEDGMENT

This work was supported by the Natural Science Foundation of China (Grant NO.61571438) and Science Foundation Research Project of Beijing, China (Grant NO.1152001).

REFERENCES

- [1] P. R. Griffiths and J. A. De Haseth, *Fourier transform infrared spectrometry*. John Wiley & Sons, 2007, vol. 171.
- [2] G. A. Pearson, "A general baseline-recognition and baseline-flattening algorithm," *Journal of Magnetic Resonance (1969)*, vol. 27, no. 2, pp. 265–272, 1977.
- [3] A. K. Atakan, W. Blass, and D. Jennings, "Elimination of baseline variations from a recorded spectrum by ultra-low frequency filtering," *Applied Spectroscopy*, vol. 34, no. 3, pp. 369–372, 1980.
- [4] I. W. Selesnick, H. L. Graber, D. S. Pfeil, and R. L. Barbour, "Simultaneous low-pass filtering and total variation denoising," *IEEE Transactions on Signal Processing*, vol. 62, no. 5, pp. 1109–1124, 2014.
- [5] F. Gan, G. Ruan, and J. Mo, "Baseline correction by improved iterative polynomial fitting with automatic threshold," *Chemometrics and Intelligent Laboratory Systems*, vol. 82, no. 1, pp. 59–65, 2006.
- [6] D. E. Brown, "Fully automated baseline correction of 1d and 2d nmr spectra using bernstein polynomials," *Journal of Magnetic Resonance*, vol. 114, no. 2, pp. 268–270, 1995.
- [7] J. J. de Rooi and P. H. Eilers, "Mixture models for baseline estimation," *Chemometrics and Intelligent Laboratory Systems*, vol. 117, pp. 56–60, 2012.
- [8] L. Andrade and E. S. Manolakos, "Signal background estimation and baseline correction algorithms for accurate dna sequencing," *Journal of Signal Processing Systems*, vol. 35, no. 3, pp. 229–243, 2003.
- [9] P. H. C. Eilers and H. F. M. Boelens, "Baseline correction with asymmetric least squares smoothing," 2005.
- [10] Z.-M. Zhang, S. Chen, and Y.-Z. Liang, "Baseline correction using adaptive iteratively reweighted penalized least squares," *Analyst*, vol. 135, no. 5, pp. 1138–1146, 2010.
- [11] S.-J. Baek, A. Park, Y.-J. Ahn, and J. Choo, "Baseline correction using asymmetrically reweighted penalized least squares smoothing," *Analyst*, vol. 140, no. 1, pp. 250–257, 2015.
- [12] Q. Han, Q. Xie, S. Peng, and B. Guo, "Simultaneous spectrum fitting and baseline correction using sparse representation," *Analyst*, vol. 142, no. 13, pp. 2460–2468, 2017.
- [13] J. Peng, S. Peng, A. Jiang, J. Wei, C. Li, and J. Tan, "Asymmetric least squares for multiple spectra baseline correction," *Analytica chimica acta*, vol. 683, no. 1, pp. 63–68, 2010.
- [14] R. Koenker and G. Bassett, "Regression quantiles," *Econometrica*, vol. 46, no. 1, pp. 33–50, 1978.
- [15] L. Komsta, "Comparison of several methods of chromatographic baseline removal with a new approach based on quantile regression," *Chromatographia*, vol. 73, pp. 721–731, 2011.
- [16] J. O. Ramsay, *Functional data analysis*. Wiley Online Library, 2006.
- [17] C. De Boor, C. De Boor, E.-U. Mathématicien, C. De Boor, and C. De Boor, *A practical guide to splines*. Springer-Verlag New York, 1978, vol. 27.
- [18] P. H. C. Eilers and B. D. Marx, "Flexible smoothing with b-splines and penalties," *Statistical Science*, vol. 11, no. 2, pp. 89–121, 1996.
- [19] D. M. Haaland and E. V. Thomas, "Partial least-squares methods for spectral analyses. 1. relation to other quantitative calibration methods and the extraction of qualitative information," *Analytical Chemistry*, vol. 60, no. 11, pp. 1193–1202, 1988.
- [20] P. Geladi, D. MacDougall, and H. Martens, "Linearization and scatter-correction for near-infrared reflectance spectra of meat," *Appl. Spectrosc.*, vol. 39, no. 3, pp. 491–500, 1985.
- [21] R. J. Barnes, M. S. Dhanoa, and S. J. Lister, "Standard normal variate transformation and de-trending of near-infrared diffuse reflectance spectra," *Appl. spectrosc.*, vol. 43, no. 5, pp. 772–777, 1989.
- [22] N. B. Gallagher, T. A. Blake, and P. L. Gassman, "Application of extended inverse scatter correction to mid-infrared reflectance spectra of soil," *Journal of chemometrics*, vol. 19, no. 5-7, pp. 271–281, 2005.
- [23] H. Martens, J. P. Nielsen, and S. B. Engelsen, "Light scattering and light absorbance separated by extended multiplicative signal correction. application to near-infrared transmission analysis of powder mixtures," *Analytical Chemistry*, vol. 75, no. 3, pp. 394–404, 2003.