

# A Simple Algorithm for Despiking Raman Spectra

Darren A. Whitaker<sup>1,2</sup> and Kevin Hayes<sup>1,2</sup>

<sup>1</sup>Pharmaceutical Manufacturing Technology Centre (PMTc), Bernal Institute,  
University of Limerick, Limerick, Ireland

<sup>2</sup>Department of Mathematics and Statistics, University of Limerick, Limerick,  
Ireland

## Abstract

Raman Spectroscopy is a widely used analytical technique, favoured when molecular specificity with minimal sample preparation is required. The majority of Raman instruments use charge-coupled device (CCD) detectors, these are susceptible to cosmic rays and as such multiple spurious spikes can occur in the measurement. These spikes are problematic as they may hinder subsequent analysis, particularly if multivariate data analysis is required. In this work we present a new algorithm to remove these spikes from spectra after acquisition. Specifically we use calculation of modified  $Z$  scores to locate spikes followed by a simple moving average filter to remove them. The algorithm is very simple and its execution is essentially instantaneous, resulting in spike-free spectra with minimal distortion of actual Raman data. The presented algorithm represents an improvement on existing spike removal methods by utilising simple, easy to understand mathematical concepts, making it ideal for experts and non-experts alike.

**Keywords:** Modified Z-Scores; Data Processing; Raman Spectra; Despiking

## 1 Introduction

Regular practitioners of Raman spectroscopy will be familiar with the problems caused by cosmic spikes. These spurious nuisance spikes typically appear at random positions and present as positive, narrow bandwidth peaks. They arise when a charge-coupled device, used as a detector in modern Raman systems, is struck by an errant high-energy particle. Predominately these are muons but may also be protons or neutrons<sup>1</sup>. Often these particles are caused by genuine cosmic rays (exotic particle produced by exploding supernovae, black holes, *etc.*) but they can also be a result of decay of radioactive atoms present in the locality of the CCD detector. The presence of these cosmic spikes hampers further multivariate data analysis. For example, they cause distortion of the principle component direction in principal component analysis<sup>2</sup>, introduce erroneous variables in multivariate curve resolution or regression techniques and can also result in misidentification in classification analysis<sup>3</sup>.

It is desirable to be able to automatically identify, reduce and/or remove these spikes from Raman spectra. This becomes even more relevant when processing large mapping datasets prevalent within pharmaceutical research<sup>4,5,6,7</sup>. Methodologies reported in the literature can be broadly separated into three categories: (i) additional acquisition based methods; (ii) methods involving hardware modification; and the category the present work falls into (iii) single-scan

32 correction via filtering or smoothing. In the first category algorithms such as robust summa-  
 33 tion<sup>8</sup> or upper-bound spectrum<sup>9</sup> methodologies take advantage of the fact the probability of the  
 34 same pixel experiencing a cosmic spike in successive measurements is very low. In the second  
 35 category methods such as analyzing the full CCD image<sup>10</sup>, division of the spectrograph slit and  
 36 image curvature correction<sup>11</sup>. Finally in the third category methods such as moving window  
 37 filtering<sup>12</sup>, spike fitting<sup>13</sup>, wavelet transforms<sup>14,15,16</sup> and median or polynomial filters<sup>17,18</sup>.

38 In this work we present a despiking algorithm based on the calculation of modified  $Z$  scores  
 39 to locate spikes and a simple moving average filter to remove the located spikes. This algorithm  
 40 is computationally efficient and inexpensive, accurate and easy to execute and program, and  
 41 should be of great utility to all users of Raman spectroscopy. Additionally the use of modified  
 42  $Z$  scores is recommended by the National Institute of Standards and Technology (NIST) as  
 43 an outlier detection methodology<sup>19</sup> and as such this method should fit easily into regulated  
 44 industries such as the pharmaceutical industry.

## 45 2 Despiking Algorithm

46 Let  $Y_1, \dots, Y_n$  represent the values of a single Raman spectrum recorded at equally spaced  
 47 wavenumbers. From this series, form the detrended differenced series  $\nabla Y_t = Y_t - Y_{t-1}$ , ( $t =$   
 48  $2, \dots, n$ ). This simple data processing step has the effect of annihilating linear and slow moving  
 49 curve linear trends, however, sharp localised spikes will be preserved.

50 Denote the median and the median absolute deviation of the differenced series by  $M =$   
 51  $\text{median}\{\nabla Y_t\}$  and  $\text{MAD} = \text{median}\{|\nabla Y_t - M|\}$  respectively, and define modified  $Z$  scores  
 52 by

$$Z_t = \frac{0.6745 \times (\nabla Y_t - M)}{\text{MAD}}.$$

53 (The multiplier 0.6745 is included to adjust for asymptotic bias that arises when MAD is cal-  
 54 culated from normally distributed data<sup>20,21</sup>.) In theory the modified  $Z$  scores can be compared  
 55 with the tabulated tail quantiles from the normal distribution. The criterion  $|Z_t| > 3.5$  was  
 56 proposed as a guideline by the American Society of Quality Control as the basis of an *outlier-*  
 57 *labeling* rule with the objective of screening large datasets for observations that are “sufficiently  
 58 suspect to merit further investigation,”<sup>22</sup>. In this paper, wavenumbers with modified  $Z$  scores  
 59 exceeding  $\tau = 6$  in magnitude were flagged as contributing to the formation of an anoma-  
 60 lous spike. In practice the scientist will have immediate control over this threshold parameter.  
 61 Lowering the value of the spike labelling threshold parameter  $\tau$  will make the algorithm more  
 62 sensitive to the presence of potential spikes.

63 Interpolated values  $\tilde{Y}_t$  are then obtained at each candidate wavenumber by calculating the  
 64 mean of its immediate neighbours, specifically  $\tilde{Y}_t = \frac{1}{w} \sum_{t-m}^{t+m} Y_t \times \mathbb{I}(Z_t < \tau)$ , where  $\mathbb{I}(u)$  is  
 65 an indicator function taking value 1 if the condition  $u$  is satisfied and 0 otherwise, and  $w =$   
 66  $\sum_{t-m}^{t+m} \mathbb{I}(Z_t < \tau)$ . This has the effect of excluding the value  $Y_t$  itself, and values of  $Y$  flagged  
 67 as contributing to the formation of a spike, from the calculation of  $\tilde{Y}_t$ . This is desirable because  
 68 in order to characterise a spike in a Raman spectrum invariably requires a sequence of 2 to 5  
 69 inflated values of  $Y$  in a row. The width of the moving average neighbourhood is controlled by  
 70 the parameter  $m$ , which was set in our applications to  $m = 5$ . Finally, in order to accomodate  
 71 the eventuality of a spike appearing at the first or last wavenumber, the values of  $Z_1$  and  $Z_n$  are  
 72 automatically set to exceed the spike labelling threshold.

### 3 Case Study

#### 3.1 Experimental

##### Sample Preparation

Theophylline (99 %, Sigma Aldrich) and Microcrystalline Cellulose (MCC101, Avicel) were blended together in 10 % w/w proportions and compacted into 12 mm diameter tablets using a single rotary punch tablet press.

##### Instrumental Set-up

Raman spectra were collected from a 12 mm tablet using a LabRAM HR Evolution (HORIBA UK Ltd., Stanmore, UK) spectrometer system. A custom spectrometer control and data acquisition script was written using the VBScripting language to enable mapping of the full tablet surface to be carried out (Figure 1). 407 individual spectra were recorded at 500  $\mu\text{m}$  intervals using a 785 nm laser line, 10 x objective, 5 s acquisition time, 100  $\mu\text{m}$  hole diameter in the range 1230 to 1330  $\text{cm}^{-1}$ .

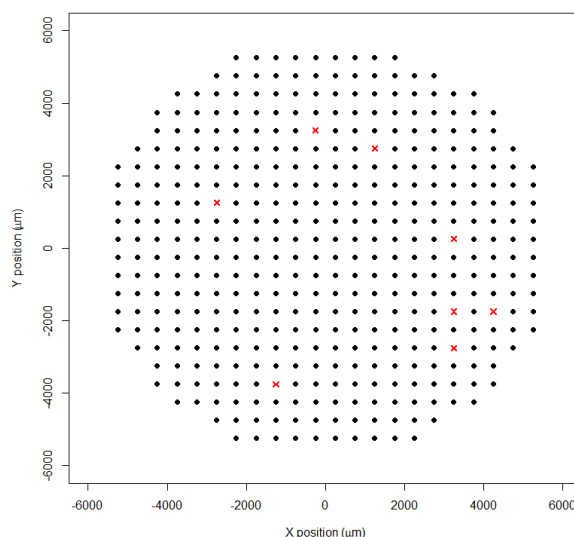


Figure 1: Sites of individual spectra over surface of 12 mm tablet (X denotes subsequently identified spike location).

##### Data Analysis

Spectral processing and analysis were performed using the R environment for statistical computing<sup>23</sup>. Custom functions for the calculation of modified  $Z$  scoring and annihilation of located spikes were developed and are included as supplementary files to this manuscript. The *hyperSpec* package<sup>24</sup> was used for easy management of spectral data within the R environment.

### 4 Results and Discussion

The Raman spectra were recorded in the wavenumber region 1230 to 1330  $\text{cm}^{-1}$ , in this region three characteristic bands of theophylline are present. The three bands centred at *ca.* 1248,

94 1286 and 1314  $\text{cm}^{-1}$  are assigned to  $\nu(\text{C-N})_{\text{sym}}$ <sup>25</sup>. An overlay of all the acquired spectra  
95 shows that cosmic spikes are present in the dataset (Figure 2a).

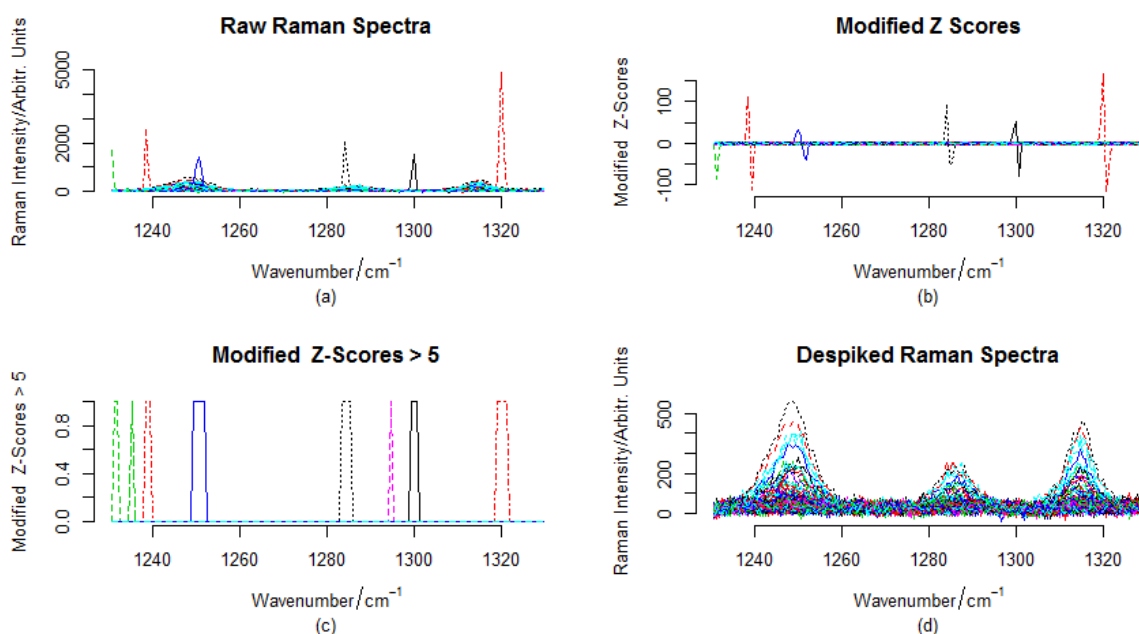


Figure 2: Raman dataset acquired on 12 mm pharmaceutical tablet and results of despiking algorithm.

96 After calculation of modified  $Z$  scores (Figure 2b) and thresholding by setting a suitable  
97 value of  $\tau$  (Figure 2c), the spikes can be removed and smoothed by applying a the moving  
98 average filter as described earlier. This results in a corrected dataset where the spikes are  
99 removed and the correct signal from the theophylline bands can be easily observed (Figure 2d).

## 100 5 Conclusion

101 We present a new algorithm based on modified  $Z$  score outlier detection for the identification  
102 and removal of cosmic spikes in Raman spectroscopic data. The algorithm was shown to be  
103 effective on a medium sized dataset acquired on a sample of pharmaceutical relevance. The  
104 algorithm is sufficiently computationally cheap to be run on almost any computer system and  
105 is also platform independent. This makes the algorithm useful for all types of Raman data  
106 analysis, including mapping measurements and real-time inline analysis.

## 107 Supplementary Information

108 The example dataset described in this paper and the despiking algorithm are provided free of  
109 charge on the publishers website.

## 110 Acknowledgements

111 This work was co-funded by the Pharmaceutical Manufacturing Technology Centre (PMTc)  
112 under Enterprise Ireland's (EI) - Technology Centres Programme & by the European Regional  
113 Development Fund (ERDF) under Ireland's European Structural and Investment Funds Pro-  
114 grammes 2014 - 2020

## References

- [1] Groom, D.. Cosmic rays and other nonsense in astronomical CCD imagers. *Experimental Astronomy* 2002;14(1):45–55. doi:10.1023/A:1026196806990.
- [2] De Groot, P.J., Postma, G.J., Melssen, W.J., Buydens, L.M.C., Deckert, V., Zenobi, R.. Application of principal component analysis to detect outliers and spectral deviations in near-field surface-enhanced Raman spectra. *Analytica Chimica Acta* 2001;446(1-2):71–83. doi:10.1016/S0003-2670(01)01267-3.
- [3] Zhang, L., Henson, M.J.. A practical algorithm to remove cosmic spikes in raman imaging data for pharmaceutical applications. *Applied Spectroscopy* 2007;61(9):1015–1020. doi:10.1366/000370207781745847.
- [4] Potter, C.B., Kollamaram, G., Zeglinski, J., Whitaker, D.A., Croker, D.M., Walker, G.M.. Investigation of polymorphic transitions of piracetam induced during wet granulation. *European Journal of Pharmaceutics and Biopharmaceutics* 2017;119:36–46. doi:10.1016/j.ejpb.2017.05.012.
- [5] Vankeirsbilck, T., Vercauteren, A., Baeyens, W., Van der Weken, G., Verpoort, F., Vergote, G., et al. Applications of Raman spectroscopy in pharmaceutical analysis. *TrAC Trends in Analytical Chemistry* 2002;21(12):869–877. URL: <http://linkinghub.elsevier.com/retrieve/pii/S0165993602012086>. doi:10.1016/S0165-9936(02)01208-6.
- [6] Gordon, K.C., McGoverin, C.M.. Raman mapping of pharmaceuticals. *International Journal of Pharmaceutics* 2011;417(1-2):151–162. URL: <http://linkinghub.elsevier.com/retrieve/pii/S0378517310009713>. doi:10.1016/j.ijpharm.2010.12.030.
- [7] Wartewig, S., Neubert, R.H.. Pharmaceutical applications of Mid-IR and Raman spectroscopy. *Advanced Drug Delivery Reviews* 2005;57(8):1144–1170. URL: <http://linkinghub.elsevier.com/retrieve/pii/S0169409X0500058X>. doi:10.1016/j.addr.2005.01.022.
- [8] Takeuchi, H., Hashimoto, S., Harada, I.. Simple and efficient method to eliminate spike noise from spectra recorded on charge-coupled device detectors. *Applied spectroscopy* 1993;47(1):129–131.
- [9] Zhang, D., Ben-Amotz, D.. Removal of cosmic spikes from hyper-spectral images using a hybrid upper-bound spectrum method. *Applied Spectroscopy* 2002;56(1):91–98. doi:10.1366/0003702021954269.
- [10] Zhao, J.. Image curvature correction and cosmic removal for high-throughput dispersive Raman spectroscopy. *Applied spectroscopy* 2003;57(11):1368–1375.
- [11] Zhang, D., Hanna, J.D., Ben-Amotz, D.. Single scan cosmic spike removal using the upper bound spectrum method. *Applied spectroscopy* 2003;57(10):1303–1305.
- [12] Katsumoto, Y., Ozaki, Y.. Practical algorithm for reducing convex spike noises on a spectrum. *Applied Spectroscopy* 2003;57(3):317–322. doi:10.1366/000370203321558236.

- 154 [13] Hill, W., Rogalla, D.. Spike-Correction of Weak Signals from Charge-Coupled Devices  
155 and Its Application to Raman Spectroscopy. *Analytical Chemistry* 1992;64(21):2575–  
156 2579. doi:10.1021/ac00045a019.
- 157 [14] Maury, A., Revilla, R.I.. Autocorrelation Analysis Combined with a  
158 Wavelet Transform Method to Detect and Remove Cosmic Rays in a Sin-  
159 gular Raman Spectrum. *Applied spectroscopy* 2015;69(8):984–92. URL:  
160 <http://asp.sagepub.com/content/69/8/984.full>. doi:10.1366/14-  
161 07834.
- 162 [15] Ehrentreich, F., Summchen, L.. Spike removal and denoising of Raman spec-  
163 tra by wavelet transform methods. *Analytical Chemistry* 2001;73(17):4364–4373.  
164 doi:10.1021/ac0013756.
- 165 [16] Tian, Y., Burch, K.S.. Automatic Spike Removal Algorithm for Raman Spectra. *Applied*  
166 *Spectroscopy* 2016;70(11):1861–1871. doi:10.1177/0003702816671065.
- 167 [17] Phillips, G.R., Harris, J.M.. Polynomial Filters for Data Sets with Outlying or  
168 Missing Observations: Application to Charge-Coupled-Device-Detected Raman Spec-  
169 tra Contaminated by Cosmic Rays. *Analytical Chemistry* 1990;62(21):2351–2357.  
170 doi:10.1021/ac00220a017.
- 171 [18] Schulze, H.G., Turner, R.F.B.. A fast, automated, polynomial-based cosmic ray spike-  
172 removal method for the high-throughput processing of raman spectra. *Applied Spec-*  
173 *troscopy* 2013;67(4):457–462. doi:10.1366/12-06798.
- 174 [19] Heckert, N.A., Filliben, J.J.. Exploratory Data Analysis. In:  
175 NIST/SEMATECH e-handbook of statistical methods; vol. 1; chap. 1.  
176 2003,URL: <http://www.itl.nist.gov/div898/handbook>.  
177 doi:papers3://publication/uuid/DE51D3F6-8B2C-4EEC-8B48-B598C13EE5F4.
- 178 [20] Tukey, J.W.. *Exploratory data analysis*; vol. 2. Reading, Mass.; 1977.
- 179 [21] Hayes, K.. Finite-sample bias-correction factors for the median absolute deviation.  
180 *Communications in Statistics: Simulation and Computation* 2014;43(10):2205–2212.  
181 doi:10.1080/03610918.2012.748913.
- 182 [22] Iglewicz, B., Hoaglin, D.. Volume 16: How to Detect and Handle Outliers. In:  
183 *The ASQC Basic References in Quality Control: Statistical Techniques*; vol. 16. ISBN  
184 9780873892476; 1993,.
- 185 [23] R Development Core Team, . *R: A Language and Environment for Statistical Com-*  
186 *puting*. R Foundation for Statistical Computing; Vienna, Austria; 2013. URL:  
187 <http://www.r-project.org>.
- 188 [24] Beleites, C., Sergo, V.. hyperSpec: a package to handle hyperspectral data sets in R;  
189 2015. URL: <http://hyperspec.r-forge.r-project.org>.
- 190 [25] Gunasekaran, S., Sankari, G., Ponnusamy, S.. Vibrational spectral investigation on  
191 xanthine and its derivatives - Theophylline, caffeine and theobromine. *Spectrochim-*  
192 *ica Acta - Part A: Molecular and Biomolecular Spectroscopy* 2005;61(1-2):117–127.  
193 doi:10.1016/j.saa.2004.03.030.