

Analyst

Accepted Manuscript



This is an *Accepted Manuscript*, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this *Accepted Manuscript* with the edited and formatted *Advance Article* as soon as it is available.

You can find more information about *Accepted Manuscripts* in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this *Accepted Manuscript* or any consequences arising from the use of any information it contains.

1
2
3 **Multivariate statistical methodologies applied in biomedical Raman spectroscopy: Assessing the**
4 **validity of partial least squares regression using simulated model datasets.**
5

6 Mark E. Keating^{1,2*}, Haq Nawaz³, Franck Bonnier^{1,4} and Hugh J. Byrne¹
7
8
9

10 **Abstract**

11
12 Raman spectroscopy is fast becoming a valuable analytical tool in a number of biomedical scenarios,
13 most notably disease diagnostics. Importantly, the technique has also shown increasing promise in
14 the assessment of drug interactions on a cellular and subcellular level, particularly when coupled
15 with multivariate statistical analysis. However, with respect to both Raman spectroscopy and the
16 associated statistical methodologies, an important consideration is the accuracy of these techniques
17 and more specifically the sensitivities which can be achieved, and ultimately the limits of detection
18 of the various methods. The purpose of this study is thus the construction of a model simulated
19 data set with the aim of testing the accuracy and sensitivity of the partial least squares regression
20 (PLSR) approach to spectral analysis. The basis of the dataset is the experimental spectral profiles of
21 a previously reported Raman spectroscopic analysis of the interaction of the cancer
22 chemotherapeutic agent cisplatin in an adenocarcinomic human alveolar basal epithelial cell- line, *in*
23 *vitro*, and is thus reflective of actual experimental data. The simulated spectroscopic data is
24 constructed by adding known perturbations which are independently linear in drug dose, as well as
25 cytological response, experimentally determined by the 3-(4,5-dimethylthiazol-2-yl)-2,5-
26 diphenyltetrazolium bromide (MTT) cytotoxicity assay. It is demonstrated that, through appropriate
27 choice of dose range, PLSR against the respective targets can differentiate between the
28 spectroscopic signatures of the direct chemical effect of the drug dose and the indirect cytological
29 effect it produces.
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54

55 **Keywords:** Raman Spectroscopy, Drug interaction studies, Partial Least Squares Regression,
56 Simulated dataset.
57
58
59
60

1
2
3 ¹FOCAS Research Institute, Dublin Institute of Technology, Kevin Street, Dublin 8, Ireland.

4
5 ²School of Physics, Dublin Institute of Technology, Kevin Street, Dublin 8 Ireland.

6
7 ³National Institute for Biotechnology and Genetic Engineering (NIBGE), P.O.Box 577, Jhang Road
8 Faisalabad, Pakistan.

9
10 ⁴Faculty of Pharmacy, EA 6295 – NM/NP, Université François-Rabelais de Tours, 60 rue du Plat
11 D'Etain, 37020 Tours Cedex 1, France

12
13
14 *E-mail: Mark.Keating@mydit.ie
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Introduction

Over the past couple of decades, vibrational spectroscopy (in particular Raman and infrared absorption) has emerged as a powerful tool for biomedical applications. The numerous studies explore applications such as disease diagnostics¹⁻⁴, cellular imaging⁵⁻⁸, the study of drug⁹⁻¹¹ and nanoparticle interactions¹²⁻¹⁴ on a cellular and sub-cellular level, to name but a few. In both modalities, the spectrum of tissue or cells contains a wealth of information and represents the combined molecular fingerprints of the ensemble of biomolecules contained in the sample. As a result, only in the simplest of cases can a valid interpretation be made by visual inspection of the spectrum. Multivariate statistical methods are thus critical in the analysis, interpretation and representation of the complex information contained within. However, given the critical nature of the outcomes of the application, whether in terms of medical diagnostics or in preliminary screening of drug efficacy and action mechanisms, it is imperative that the combination of spectroscopic techniques and multivariate analysis are rigorously and quantifiably validated. Such validation can also establish realistic limits to what is often purported as a high content screening methodology. To this aim, the use of simulated datasets based on experimental studies can play a crucial role^{14,15}.

A multitude of multivariate analytical methods exists, each of which aims to simplify complex biospectroscopic information and provide a tool with which to draw conclusions about the state of the sample. These include Principal Component Analysis (PCA), Linear Discriminant Analysis (LDA), Vertex Component Analysis (VCA), Spectral Cross Correlation Analysis (SCCA), K-means Clustering Analysis (KMCA), Hierarchical Cluster Analysis (HCA) to name but a few. Importantly, there also exists a number of variants of these methods which differ slightly and can give, in some instances, different answers^{14,16,17}.

Recently, regression modelling (e.g. Partial Least Squares Regression, PLSR) has seen a number of biomedical uses in both Raman and IR spectroscopies. The core idea of using this method is to investigate the spectral variability as a function of a systematic conditional change such as

1
2
3 radiation dose¹⁸ or viral infection¹⁹. PLSR can be employed to construct predictive models for
4
5 spectral response as a function of the target variable. Therefore, an unknown dose or degree of
6
7 infection can be determined from its spectrum, having obvious potential clinical applications.
8
9 Furthermore, feature selection techniques such as PLSR o-efficients, Jack-Knifing (JK) and genetic
10
11 algorithms, amongst others²⁰, can be employed to identify the most statistically relevant spectral
12
13 changes, such that the biological mechanisms underlying the spectral changes can explored and
14
15 understood. Importantly, there are many variants of the PLSR algorithm and, in some instances,
16
17 hybrid methods which use a combination of two statistical tools in order to extract relevant chemical
18
19 information have been employed. Although these methods have been applied to a wide range of
20
21 studies, the details are beyond the scope of this paper although good examples can be found in
22
23 literature^{10,11,18,21-24}
24
25
26

27
28 The potential of Raman spectroscopic microscopy for initial screening of chemotherapeutic
29
30 efficacy and mechanism of action has been demonstrated by Nawaz *et al.*^{10,11,23}. Taking the
31
32 interaction of cisplatin with the human lung adenocarcinoma cell line, A549, *in vitro*, as an example,
33
34 PLSR of Raman spectroscopic datasets was reported to identify and differentiate the direct effects of
35
36 cisplatin on the cellular biochemistry as a function of drug concentration (dose) and the resultant
37
38 toxicological response as measured by the MTT cytotoxicity assay. This simultaneously provides a
39
40 parallel gold standard technique to compare to the spectroscopic endpoint as well as range finding
41
42 for the initial dose response curve i.e. establishing values of Inhibitory Concentrations (IC) etc. In an
43
44 operational model of pharmacological agonism, the former is a linear process, where as the latter
45
46 results in the more complex sigmoidal response of cell populations to drug exposure²⁵. PLSR against
47
48 the drug concentration returned changes in the Raman peaks associated with both conformational
49
50 and chemical changes in DNA, while changes to the lipid and protein distributions were dominant
51
52 when the data was regressed against the cytotoxicological end point, indicating the biochemical
53
54 changes associated with the resultant cytological response to the interaction with cisplatin. The
55
56 statistic relevance of the results were confirmed using the JK approach.
57
58
59
60

1
2
3 The potential to differentiate the direct chemical effects from the subsequent cytological
4 responses opens the way to the use of the techniques to visualise and interpret the mode of action
5 of chemotherapeutic agents intracellularly and to quantify the efficacy to produce the desired
6 cellular response in a single truly label free measurement. The emergence of ever higher throughput
7 spectrometers would enable realtime and time resolved visualisation of the respective processes as
8 they evolve. Notably, however, while the studies of Nawaz *et al.* show great promise towards this
9 end, the technique is as yet unvalidated. The expected changes in the spectra with concentration
10 and toxicological endpoint are inferred, based on prior knowledge about the biological action of
11 cisplatin in the model *in vitro* system. This leads to a difficulty when trying to confirm the validity of
12 the method or compare two different methods to quantitatively assess the sensitivity, accuracy and
13 specificity of the technique.
14
15
16
17
18
19
20
21
22
23
24
25
26

27 Here, we aim to validate the application of these methodologies using simulated datasets
28 based on the previously published experimental results of Nawaz *et al.*. In particular, we aim to test
29 the ability of PLSR to model and thus extract spectroscopic variations (based on the regression co-
30 efficients) which vary systematically as a function of different targets. Thus, the study will confirm
31 whether the method is capable of extracting and differentiating spectroscopic features which differ
32 based on linear or non-linear changes of the targets. Additionally, the accuracy or fidelity of the
33 method in extracting systematically varied features will be explored as the spectral perturbations
34 introduced decrease in magnitude, exploring the sensitivity of the method. Thus, the overarching
35 aim is to establish the validity of the algorithms applied to Raman spectral datasets containing
36 changes pertaining to the direct and indirect effects of the anti-cancer drug cisplatin *in vitro*. For the
37 purposes of this study, we propose the use of a modelled simulated dataset. The dataset is
38 constructed based on experimental observations, but the systematic spectral variation that is
39 introduced is known precisely and thus an exact and complete assessment of the method can be
40 carried out.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Experimental

Experimental results were obtained as described in previous publications by Nawaz *et al*^{10,11} which investigated Raman spectroscopy as a tool to study cisplatin-cellular interactions *in vitro*. The experimental methods are described in detail in the publications, but are summarised in brief as follows.

Human lung adenocarcinoma (A549) cells were routinely cultured at 37 °C, 5 % CO₂ in DMEM F12 supplemented with 10% FBS, 1% pen/strep and 2mM l-glutamine. Cells were cultured until 70-80% confluency and plated on quartz substrates for Raman spectroscopy . A standard MTT assay, using a concentration range of 0.05µM – 50 µM, was used to assess the toxicity of cisplatin to provide a comparison to Raman spectroscopy. This was carried out in standard 96 well plates and experiments were all completed in triplicate. This range resulted in a sigmoidal variation in cell culture viability over the range ~90% to ~20%, from which the Inhibitory Concentration (IC₅₀) of cisplatin in A549 cells *in vitro* was determined to be 1.2 ± 0.2 µM.

Cisplatin, at varying concentrations in the range 0.05 µM - 50µM, was added to cells and Raman microscopic measurements of cells exposed to each dose, including unexposed control, were acquired at a source wavelength of 785nm for both nuclear¹⁰and cytoplasmic regions¹¹ . In both cases, multiple spectra were recorded from a total of 60 cells at each exposure level. The PLSR approach was used to model the spectroscopic data as well as to select and distinguish the relevant features indicative of the chemical effects of cisplatin and the cellular response to cisplatin via a regression against dose and the MTT cytotoxicity endpoint respectively. By examination of the regression co-efficients, it was possible to discern the major features responsible for model construction.

In this work, these experimental spectral datasets are employed to construct semi-realistic simulated data to probe the reliability, sensitivity and quantitative nature of these methods when

1
2
3 applied to drug-interaction studies. More details of the experimental set up can be found in Nawaz
4
5 *et al.*^{10,11}
6
7

8 **Partial Least Squares Regression**

9
10 PLSR is a multivariate statistical method which aims to establish a model that relates the variations
11
12 of the spectral data to a series of relevant targets. The spectral data (X matrix) is thus related to the
13
14 targets (Y matrix) according to the linear equation $Y = XB + E$, where B is a matrix of regression
15
16 coefficients and E is a matrix of residuals. The PLSR algorithms used in this study have been
17
18 previously published elsewhere^{10,11,18,22} and are based on scripts written in house using Matlab 7.2
19
20 (The Mathworks Inc.). The algorithm allows for the construction of a regression model which can be
21
22 used to predict the outcome in a number of different situations. In this case, the examples used are
23
24 concentration and MTT response, and therefore the algorithm can be used to predict for example
25
26 the toxicological response of a particular drug dose.
27
28

29
30 Latent variables (LV's) in PLSR modelling are a series of underlying variables which aim to
31
32 describe the behaviour of the modelled system. The exact number of latent variables which are
33
34 necessary to build an entirely accurate model is not known *a priori*. However, it is one of the goals of
35
36 PLSR models to accurately predict the number necessary to build a robust and accurate model²⁶.
37
38 Predicting the number of LVs which will build an accurate model is usually achieved during the cross
39
40 validation step, typically using the root mean squared error of cross validation (RMSECV) as a metric
41
42 for latent variable selection.
43
44
45
46
47

48 **Spectral Constructs**

49
50 Spectral constructs were generated for the purpose of imparting a known perturbation to the
51
52 dataset which could be systematically varied to evaluate the capability of the PLSR modeling to
53
54 accurately predict and extract spectral variations correlated to a known external variable, in this
55
56 case, drug dose and the resultant cytological changes. Using the original datasets of Nawaz *et al.*,
57
58
59
60

1
2
3 derived from the nuclear and cytoplasmic regions, specific spectral changes were identified in the
4
5 mean difference spectra of a 3 μ M exposed cell population versus the unexposed control (Figure 3,
6
7 of reference 10, Figure 4 of reference 11). In this way, spectral constructs were generated from the
8
9 changes in the spectra of the nuclear region, including increases in the characteristic A form of DNA
10
11 peak at 807 cm⁻¹ and the B form peak at 833 cm⁻¹ and a change in the C-H deformation at 1449 cm⁻¹
12
13 (Figure 1A) and in the cytoplasmic region, containing the following peak changes or shifts; a change
14
15 in the amide 1 band at ~1661 cm⁻¹, a decrease in the C-C stretch intensity at ~939 cm⁻¹ and an
16
17 increase in the tryptophan peak at 731 cm⁻¹ (Figure 1B). The relative intensities of the peaks in each
18
19 construct were derived from the experimental difference spectra at a cisplatin exposure dose of
20
21 3 μ M¹⁰ and were normalised for concentration (Figure 1A) and a loss of viability at that concentration
22
23 of 0.52¹⁰ (Figure 1B). Different weightings of these spectral constructs (termed hereafter the
24
25 Concentration and Viability construct respectively) were then added to a control dataset as
26
27 described in the following section.
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

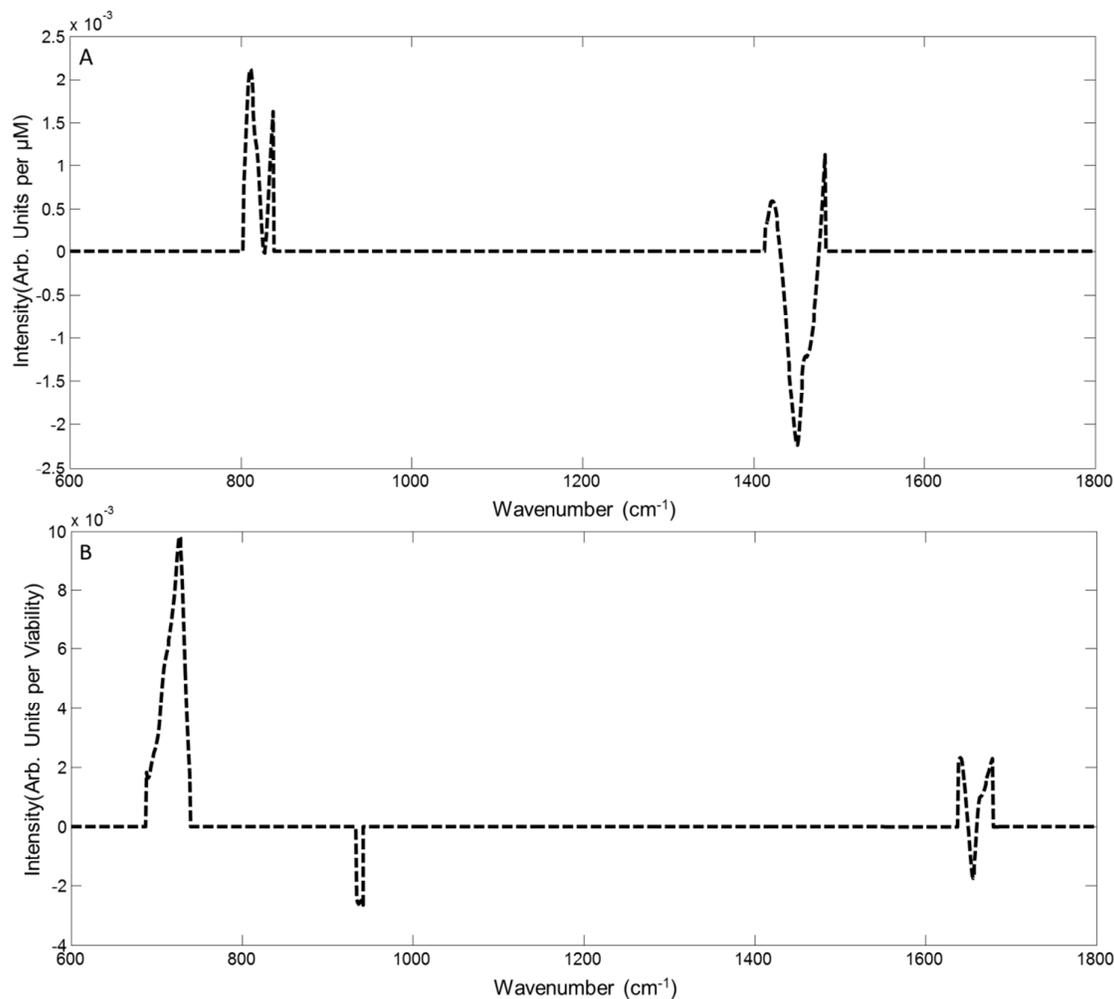
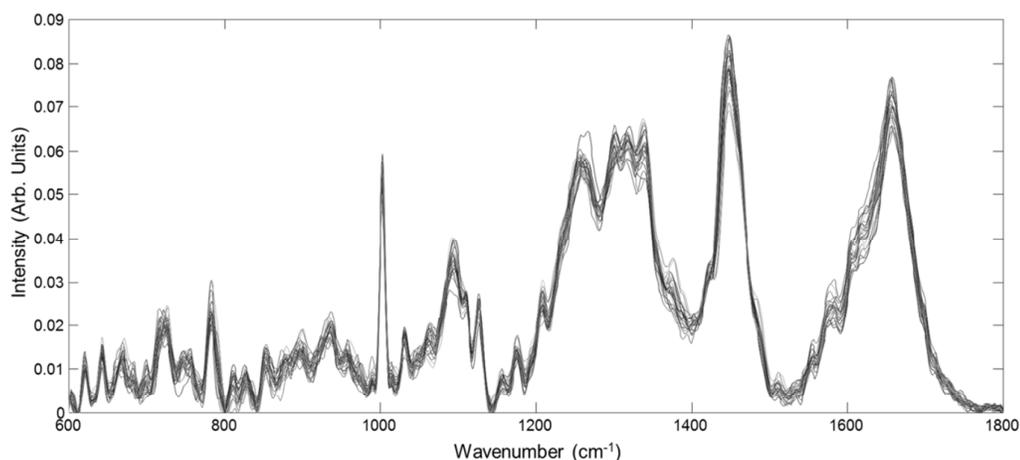


Figure 1: Spectral Constructs based on the normalised difference spectra between control and exposed nucleus (A)¹⁰, and cytoplasm¹¹ (B). Selected Raman peaks were used to avoid over complexity in the simulated data; (A) the A form peak of DNA at 807 cm^{-1} and the B form peak at 833 cm^{-1} and the C-H deformation at 1449 cm^{-1} (B) the amide 1 band at $\sim 1661 \text{ cm}^{-1}$, the C-C stretch intensity at $\sim 939 \text{ cm}^{-1}$ and the tryptophan peak at 731 cm^{-1} .

Simulated data

Simulated datasets were generated in the following manner. A control dataset containing 25 spectra acquired from the nucleus of non-cisplatin exposed (control) cells was selected from Nawaz *et al*¹⁰ (Figure 2). Notably, this real experimental dataset contains instrumental noise and sample variability. To this dataset, weighted contributions of the Concentration construct shown in Figure 1A, based on the experimentally observed difference spectra of the nuclear region, were added,

1
2
3 over the Lethal Concentration range 0.05 μM - 50 μM used in the original study, which includes the
4
5 IC_{50} , based on a direct weighting of the spectral construct by the range of concentrations (Table 1).
6
7 Initially, only the concentration dependent weighted constructs were added to the control, to
8
9 produce Dataset 1.
10



11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 2: Control dataset taken from Nawaz *et al.*¹⁰; 25 control spectra taken from the nucleus of cells not exposed to cisplatin. Spectra have been baseline corrected and vector normalised. The inherent spectral variability in the data is representative of real experimental conditions. These spectra were then used in the construction of 3 simulated datasets, each containing 8 different dose/viability points with systematically introduced variation of the spectral constructs shown in figure 1.

As the MTT assay is expressed in viability compared to control (0.845 being maximum (V_{max}) and 0.135 being minimum values of fit to the experimentally observed viability over the concentration range¹⁰), the spectral construct of Figure 1B, derived from the experimentally observed differences in the cytoplasmic region, was similarly weighted by the ($V_{\text{max}} - \text{MTT}$) endpoints in Table 1 and also added to Dataset 1. Each spectral construct was therefore added following a linear trend based on concentration (Figure 1A) plus a linear trend based on MTT response (Figure 1B). The MTT endpoint data are, however, nonlinearly related to the concentration, in a sigmoidal fashion typical of cytotoxic responses, as shown in Nawaz *et al.*^{10,11}. The resultant dataset therefore contains 25 spectra for each of 8 dose points (including control) which incorporate

spectral variations, systematically dependent on both the exposure dose and the measured cytological response. For simplicity, this is referred to as Dataset 2.

It is noted that the spectral construct of Figure 1B is derived from exposure dose dependent, experimentally observed, spectral changes in the cytoplasmic region. No direct biological significance is inferred by the weighted addition of this spectral construct to the dataset derived from the nuclear regions. However, the addition serves to provide an independently variable perturbation to the dataset, which may serve to mimic a cytological effect of the direct action of the drug in the nucleus.

To probe the sensitivity of the methodology, the experimental range for cisplatin (Lethal Concentration, in table 1) has been extended (Sub lethal Concentration in table 1) to represent non-lethal doses of the drug. The MTT values have also been extrapolated according to the original fit of the Hill equation¹⁰ to reflect these changes in concentration (Sub-lethal MTT in table 1). The corresponding simulated dataset will be referred to as Dataset 3. A dataset was also constructed which consisted solely of control spectra. This Control dataset did not contain any systematically introduced spectral variations and was used to establish a baseline regression endpoint for both Lethal Concentration and Lethal MTT.

| <i>Lethal Concentration</i> | <i>Sub-lethal Concentration</i> | <i>Lethal MTT</i> | <i>Sub-lethal MTT</i> |
|------------------------------------|--|--------------------------|------------------------------|
| <i>0.05</i> | <i>0.0005</i> | <i>0</i> | 0.000001 |
| <i>0.5</i> | <i>0.005</i> | <i>0.15</i> | 0.000001 |
| <i>1</i> | <i>0.01</i> | <i>0.35</i> | 0.000001 |
| <i>3</i> | <i>0.03</i> | <i>0.52</i> | 0.00001 |
| <i>5</i> | <i>0.05</i> | <i>0.55</i> | 0.0001 |
| <i>10</i> | <i>0.1</i> | <i>0.65</i> | 0.001 |
| <i>50</i> | <i>0.5</i> | <i>0.66</i> | 0.01 |

1
2
3 *Table 1: The weightings of the spectral constructs added to the control data. The Lethal*
4 *Concentration and Lethal MTT ranges are derived from the actual experiment data of references*^{10,11}.
5 *Lethal MTT represents the values obtained when the experimental MTT value is subtracted from V_{max} .*
6 *The Sublethal Concentrations extend the concentration range and are representative of sub-lethal*
7 *doses of cisplatin, for which sub-lethal MTT values are derived from the extrapolated fit of the Hill*
8 *equation in Reference 10.*
9

14 Results

17 Concentration Simulated data

18
19
20 The PLSR method aims to establish a model that relates the variations of the spectral data to a series
21 of relevant targets. In this case, the spectral data is a series of simulated datasets which are based on
22 known introduced perturbations based on cisplatin-cellular interactions as described in the previous
23 sections.
24
25
26
27

28
29 Regression of Dataset 1 against the Lethal Concentration range (table 1) yielded the model
30 shown in figure 3. The data were split, 60:40, to create calibration and test sets to build the model.
31
32 60% of the data was used to calibrate the model and 40% of the data was then used to assess the
33 performance of the model in predicting the expected target with unseen data. Leave-one out cross
34 validation with the calibration set was used to determine the optimal model complexity for use in
35 testing (Meade et al., 2010)²⁷. This process was performed with randomization of the data matrix
36 and splitting of the data to prevent data bias (Varmuza and Filzmoser, 2009)²⁸. Control of over fitting
37 was achieved using a procedure previously described by Martens and Naes²⁹. The procedure involves
38 selection of the optimal number of latent variables (LV) to retain within the PLSR model via cross-
39 validation with the calibration data set. The optimal number of LV's was then selected on the basis
40 of the number which provided the lowest root mean squared error after cross validation. This is
41 illustrated in Supplementary Material figure S1A and B, which show plots of the RMSECV and RMSEP
42 for the first 10 LV's for the regression of Dataset 1 against Lethal Concentration 1. The values for
43 RMSECV and RMSEP approach zero in an asymptotic fashion, and as there is no significant further
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

decrease after 10 LVs, 10 was chosen as the optimum number. and thus the optimum number of LV's was selected as 10. The calibration and test set had RMSEC=0.49673, RMSEP=0.52389 and R^2 values of 0.99948 and 0.99953 respectively, indicating a good linear fit of the model.

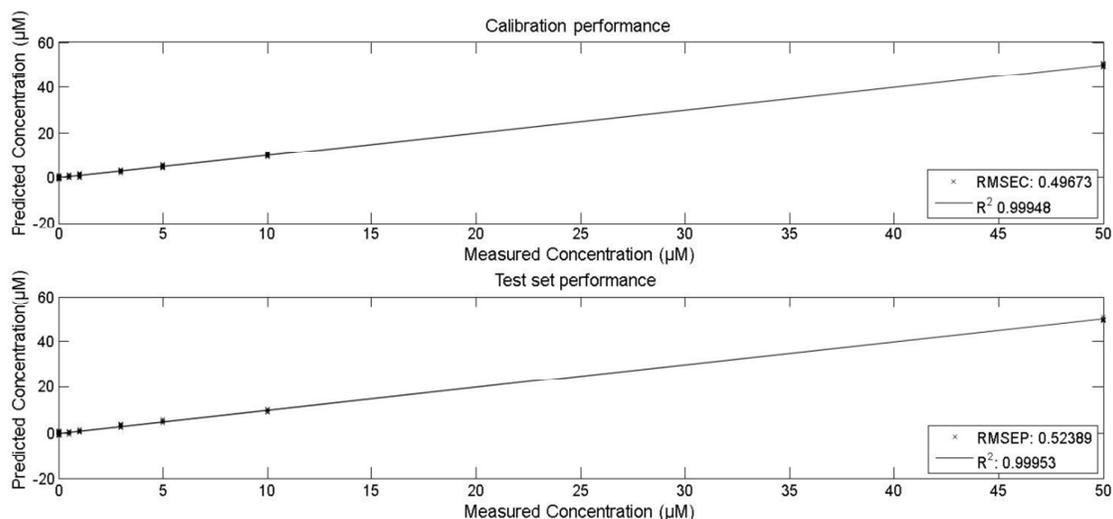


Figure 3. PLSR modelling against Lethal Concentration for Dataset 1. Top panel shows the calibration performance and test dataset (RMSEC 0.49673, R^2 0.99948). Bottom panel shows the performance of the model for the test dataset (RMSEP 0.52389, R^2 0.99953). Data was split in a ratio of 60:40 calibration and test respectively.

As the regression co-efficients (RC) are descriptors of the spectral features which are used to build the model, we also aimed to assess the accuracy with which the algorithm can faithfully extract the known spectral perturbations introduced in the dataset. For regression of Dataset 1 against Lethal Concentration, we expect that the spectrum of the RC will be comprised of the Concentration construct which has been added based on the Lethal Concentration range (Figure 1A).

In figure 4, a direct comparison between the RC of regression of Dataset 1 against the Lethal Concentration range and the concentration spectral construct is shown. The spectrum of the RC is dominated by the peaks of the systematically added spectral construct, at 807cm^{-1} , 833cm^{-1} , which

1
2
3 correspond to A and B form DNA¹⁰ and the C-H deformation at 1449cm⁻¹ (solid line figure 4 bottom
4 panel). This verifies that the simulated changes are the major contributors to the PLSR model
5 construction.
6
7
8
9

10
11 However, it should be noted that the RC spectrum in figure 4 also contains other peaks
12 which are not present in the spectral construct and so should not show a systematic variation with
13 concentration. By regression of just the control data (with no spectral perturbations) against the Y
14 target (Lethal Concentration) it was possible to establish a Control RC, as shown by the dotted line
15 (bottom panel) in figure 4 (offset and multiplied by a factor of 10 for clarity). The control RC
16 spectrum shows a high degree of similarity with the original cellular spectra (Figure 2) and thus
17 derives from the inherent variability in the experimental measurement. Close examination of the RC
18 for the Dataset 1 regression reveals that some of the peaks in the Control RC are also present.
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

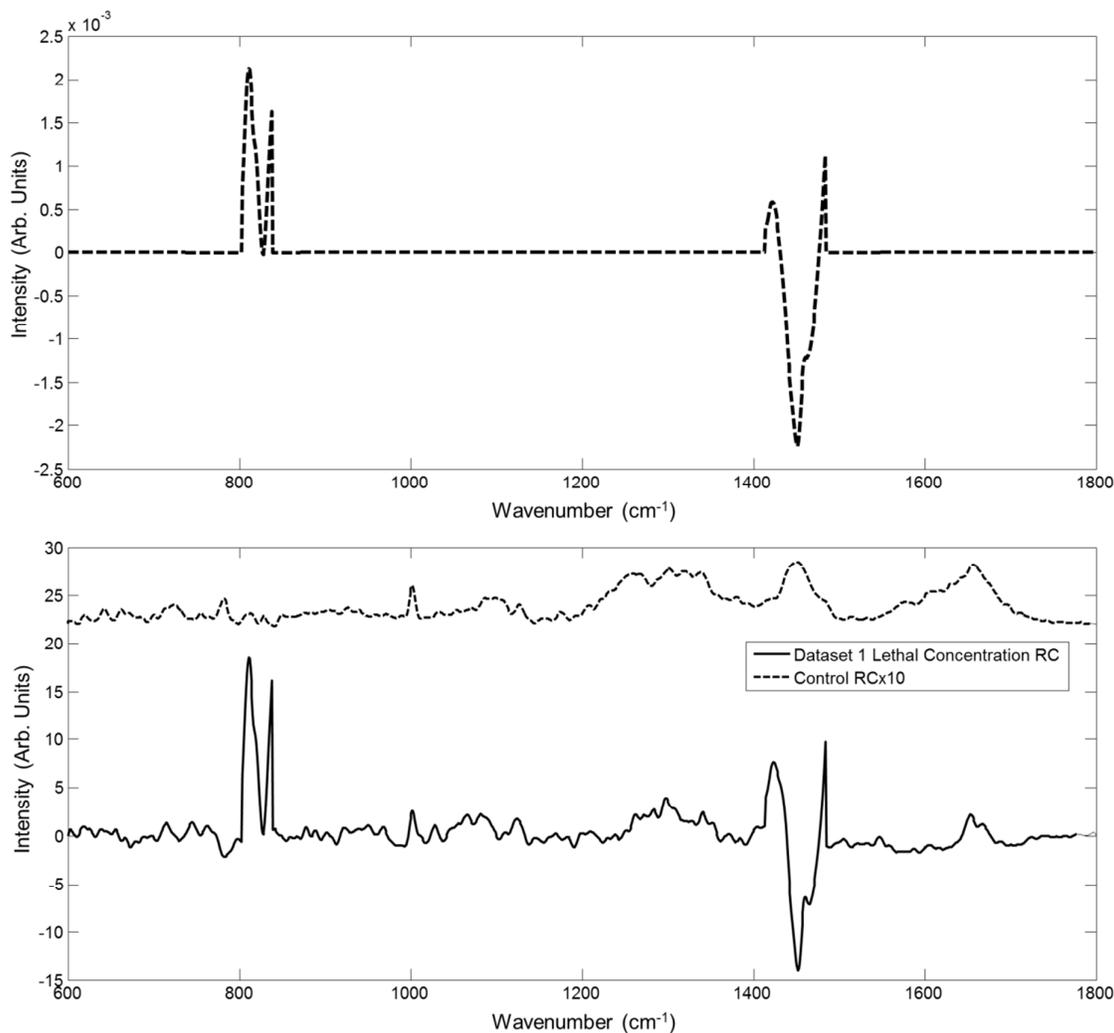


Figure 4: Plot of the regression co-efficients following PLSR of Dataset 1 against Lethal Concentration.

The Concentration construct (dashed line) is shown in the top panel for comparison with the RC's in the bottom panel. The solid line (bottom panel) shows the regression co-efficient following regression of Dataset 1 against Lethal Concentration. The dotted line shows a plot of the regression co-efficient following regression of a dataset consisting of just control spectra against Lethal Concentration, in effect showing the baseline regression co-efficient when no introduced spectral perturbation (not including sample/instrumental variations) is present. The Control RC has been offset and multiplied by a factor of 10 for clarity.

1
2
3 The PLSR modelling process was repeated for Dataset 2, which included the combined perturbations
4 of the Concentration construct of Figure 1A, linearly weighted according to Lethal Concentration of
5 Table 1, and the MTT Construct of Figure 1B, linearly weighted according to Lethal MTT of Table 1. A
6 similar performance of model calibration and test were achieved, with RMSEC=0.4981,
7 RMSEP=0.53505 and R^2 values of 0.99947 and 0.99952 respectively, again indicating a good linear fit
8 of the model (Figure S2). The spectrum of RC again faithfully reproduced the Concentration
9 Construct of Figure 1A, on a background which matches well the Control RC spectrum (Figure S3).
10
11
12
13
14
15
16
17
18
19
20
21

22 **MTT Simulated Data**

23
24 Dataset 2 also contains systematic perturbations which have been weighted according to the
25 viability as measured using the MTT assay, and it is of critical interest whether these spectral
26 variations can be independently extracted using PLSR, as suggested by Nawaz et al.¹⁰. Regression of
27 Dataset 2 against Lethal MTT (table 1) yielded the model shown in figure 5. As for the concentration
28 dependent model, the data are split according to 60% calibration and 40% test data. The calibration
29 and test set had RMSEC=0.10158, RMSEP=0.12087 and R^2 values of 0.91928 and 0.89793
30 respectively. Based on these values, it can be seen that, while the model has fitted the data, it does
31 not provide as good prediction as shown for concentration (figure 3). This is also reflected by the
32 lower R^2 values, considering that the accuracy of the linear fit is measured by how close the value is
33 to 1. A possible explanation for this is the lower magnitude and range of weightings of spectral
34 construct added corresponding to the MTT response (Table 1, Lethal MTT).
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

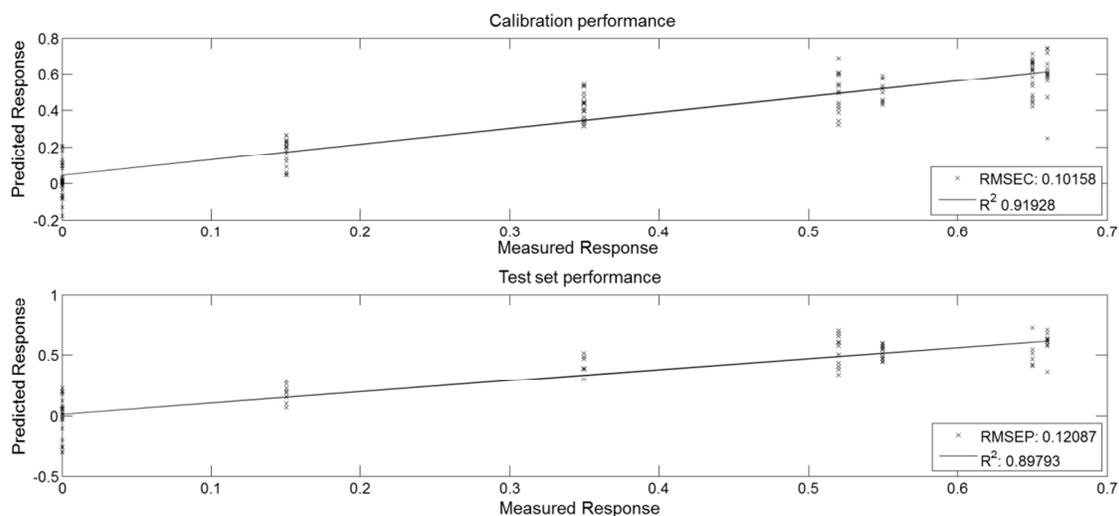


Figure 5: PLSR modelling of Dataset 2 against the Lethal MTT target. Top panel shows the calibration performance and test dataset (RMSEC 0.10158, R^2 0.91928). Bottom panel shows the performance of the model for the test dataset (RMSEP 0.12087, R^2 0.89793). Data has been split in a ratio of 60:40 calibration and test respectively.

Inspection of the MTT RC in Figure 6 shows that the peaks of the systematically added Viability construct (Figure 6, dashed line, top panel), the amide 1 band at $\sim 1661\text{ cm}^{-1}$, the C-C stretch intensity at $\sim 939\text{ cm}^{-1}$ and the tryptophan peak at 731 cm^{-1} , are faithfully reproduced and dominate the MTT RC (Figure 6, solid line, bottom panel).

The baseline sensitivity is evaluated by regressing the control dataset against the Lethal MTT target, yielding the Control RC of Figure 6 (bottom panel, dotted line). The resultant RC spectrum has been offset and multiplied by a factor of 10, for clarity. As in the case for regression against Lethal Concentration targets, the Control RC resembles the cellular spectra of figure 2, indicating that the baseline variation is limited by the variations in the original spectral measurement.

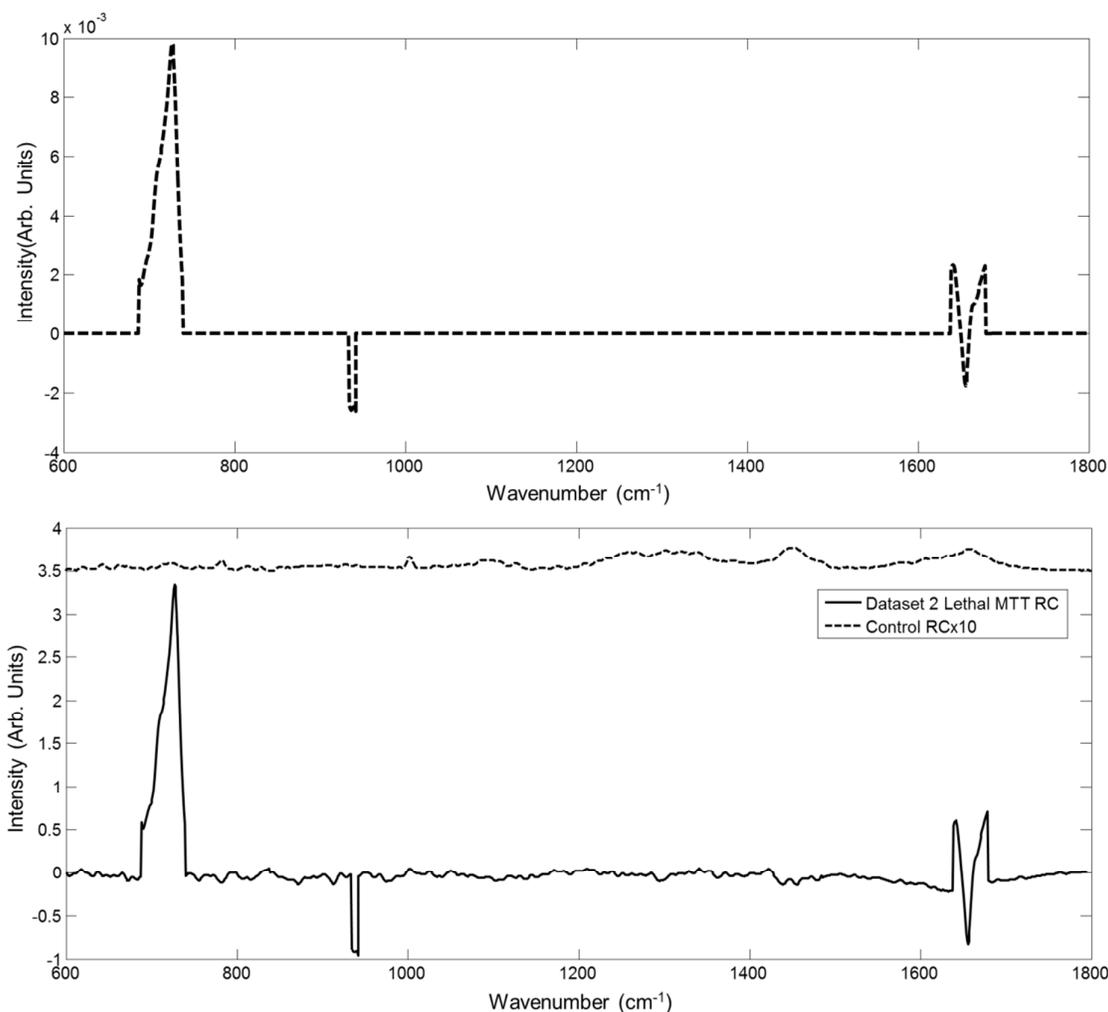


Figure 6: Plot of the regression co-efficients following PLSR modelling against MTT response. The Viability construct (dashed line) is shown in the top panel for comparison with the RC's in the bottom panel. The solid line shows the regression co-efficient following regression against Lethal MTT and Dataset 2 (bottom panel). The dotted line (bottom panel) shows a plot of the regression co-efficient following regression of a dataset consisting of just control spectra against Lethal MTT, in effect showing the baseline regression co-efficient when no introduced spectral perturbation (not including sample/instrumental variations) is present. The Control RC is offset and multiplied by a factor of 10 for clarity.

Quantative evaluation of regression co-efficients

In an attempt to evaluate the quantitative nature of the regression co-efficients, a method was devised which looked at varying the number of datapoints used to build the PLSR model. For the analysis of the spectral variations of Dataset 1, based on variations of the Concentration construct of figure 1A weighted according to Lethal Concentration (Table 1), multiple regressions were conducted (models not shown). Each model was constructed by increasing the number of data points, C+1 being the first data set used, consisting of the control dataset (Fig 2) and the 0.05 μM datapoint of the Lethal Concentration range (Table 1). The data set was then successively extended by 1 datapoint, such that C+2 consists of control, 0.05 μM and 0.5 μM , and so on, until all data points in the Lethal Concentration were included.

For all models, the spectrum of the RC displayed a combination of the Concentration construct of Figure 1A and the Control RC of Figure 4, and, as expected, regression over the full range reproduced the RC spectrum of Figure 4. Notably, as shown in Figure 7, the peaks of the Concentration construct increase linearly as the range of the regression is increased and reach a saturation value above $\sim \text{C}+4$. Extension of the model to 1000 μM results in no further significant increase of these maximum peak intensities (data not shown). The A-form DNA peak at 807 cm^{-1} reaches a maximum value of 18.46. Although this does not quantitatively equate to the corresponding peak value of the Control construct of Figure 1A, the relative magnitudes of the respective peaks is consistent with those of the original Concentration construct, and notably the relative contribution of the Control RC is reduced with increasing range.

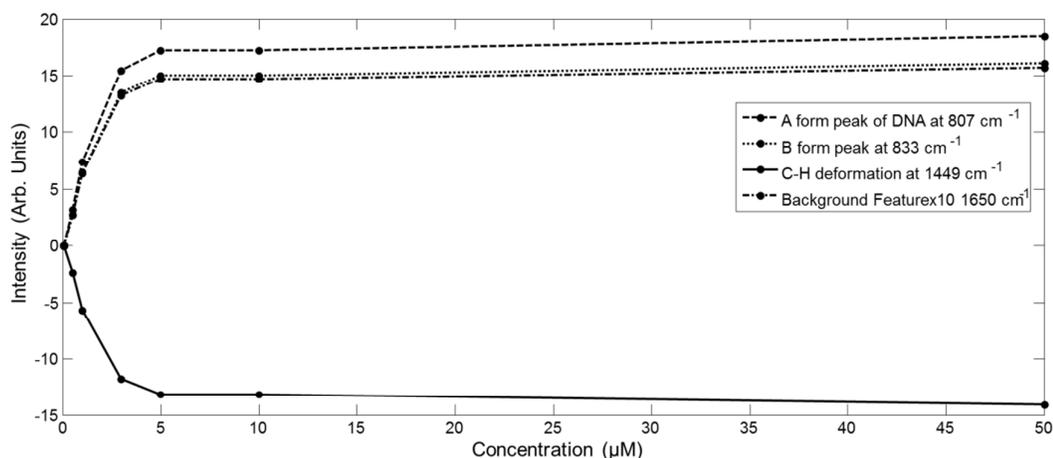


Figure 7: Evolution of the peaks of Construction construct in PLSR models of increasing range for Dataset 1.

A similar analysis was conducted for the PLSR of Dataset 2 against the Lethal Concentration range. Figure 8 shows a plot of the extracted RCs for all successive regressions. As expected, C+7 reproduces the Lethal Concentration RC of Figure 4, and extracts the expected introduced spectral construct (Figure 1 A). However, notably for all other regressions, C+1 to C+6, the presence of peaks which are not explicitly dependent on Lethal Concentration are observed. In addition to those of the Control RC, peaks of the MTT construct (Figure 1B) are evident in the RCs of the regressions over the incomplete concentration range. A similar phenomenon can be seen in the equivalent sequential modelling of the MTT data of Dataset 2 (Figure S4 and S5).

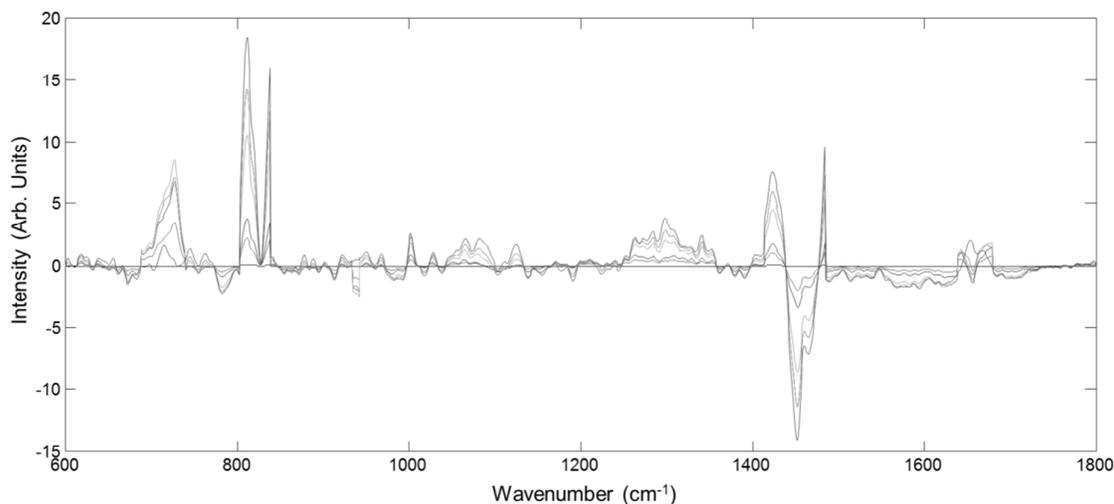


Figure 8. A plot of regression co-efficients following multiple regression against concentration with increasing data points. I.e. C+1 represents a dataset consisting of the control dataset and the data point at 0.05 μM . This then increases C+n until all data points in the dataset have been evaluated.

Figure 9 shows a plot of selected RC peak intensities associated with the spectral construct relating to concentration following successive rounds of regression as described above, namely the A form peak of DNA at 807 cm^{-1} and the B form peak at 833 cm^{-1} , which are associated with the physical changes associated with cisplatin-cellular interaction¹¹. In fact the evolution of the peaks is observed to be identical to that observed for Dataset 1, shown in Figure 7, and although the plot of Figure 9 is in a linear/logarithmic format, it can be seen that the predicted relative intensities again increase linearly initially, before reaching a point of saturation at, or above, the dataset C+4, and further addition of datapoints makes no difference (data not shown) to the quantitative prediction of the features.

Also shown in Figure 9 is the dependence of the peak of the Viability construct at 731 cm^{-1} , (for example) which “bleeds through” in the regression of Dataset 2 against the incomplete concentration range. This bleed through occurs for all peaks of the MTT Construct. The contribution

of the peaks of the Viability Construct follows a trend of the derivative of the viability curve, indicating that it is the rate of change of the contributed spectral variations which governs the contribution to the RC. Notably, when the full Lethal Concentration range is included in the model, at the extremes of which the change in viability has reduced to the minimum value, the bleed through of the MTT construct is minimal, and the Concentration Construct of Figure 1A is faithfully extracted, albeit with an underlying background as a result of the inherent spectral variability.

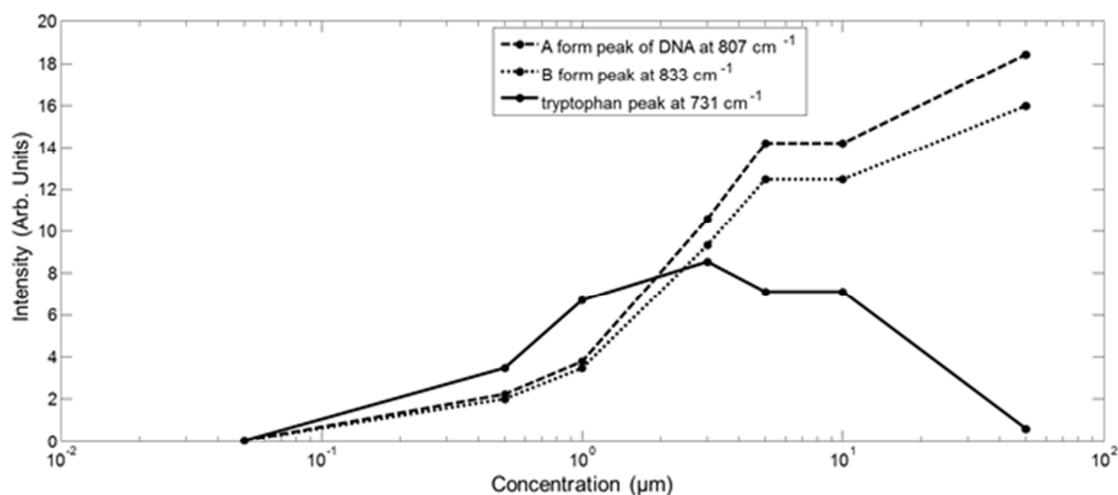


Figure 9. Plot of peak intensities vs. concentration of regression co-efficients for the A form peak of DNA at 807 cm⁻¹ and the B form peak at 833 cm⁻¹ of the Concentration Construct (Figure 1A). Also plotted is the contribution of the tryptophan peak at 731 cm⁻¹, a key feature of the Viability Construct (Figure 1B)

A similar PLSRA of the contributions of the Viability construct to Dataset 2 reveals similar bleed through and more complex evolution of the features contributing to the spectrum of the RC (Supplementary Material Figures S4 and S5). The bleed through of the features of the spectral constructs shown in Figures 8 and 9 is a clear demonstration that it is not trivial to independently extract the contributions of the two constructs over the lethal concentration range, as speculated by Nawaz *et al.*¹⁰. However, over concentration ranges in which the viability does not change

1
2
3 significantly, the bleed through is minimal, and the concentration dependent spectral changes can
4
5 be independently extracted. Thus, it should be possible to determine the direct chemical
6
7 interactions of an external agent in the sublethal range.
8
9

10
11 Figure S6 shows the calibration and test performance of the PLSR of Dataset 3 versus the Sublethal
12
13 concentration range of Table 1. The model yields RMSEC and RMSEP values of 0.143 and 0.19392,
14
15 respectively, with R^2 values of 0.38916 and -0.24063, accuracies considerably less than those of the
16
17 equivalent model in the Lethal Concentration range. Notably, the RC spectrum is a faithful extraction
18
19 of the pure Concentration construct of Figure 1 A, as shown in Figure S7. Little or no bleed through
20
21 of features associated with the Viability construct is apparent (although still present in minimal
22
23 quantities) although this is not surprising as, with little or no change in viability, the contributions of
24
25 the Viability construct to Dataset 3 are minimal.
26
27
28
29
30
31
32

33 Discussion

34
35
36
37 Given the drive for a reduction in the use of animal models for evaluating toxicity, screening of drugs
38
39 and even cosmetics, due to regulatory developments in both the EU and US (EU Directive-
40
41 2010/63/EU and US Public Law 106-545, 2010, 106th Congress)³⁰⁻³² generally based on the 3 R's of
42
43 Russell and Burch³⁰ to replace, reduce and refine the use of animals used for scientific purposes,
44
45 there is increased emphasis on the development of reliable and rapid *in vitro* screening
46
47 methodologies. This includes more representative culture models which better mimic the *in vivo*
48
49 environment as well as more rapid, cost efficient, high content, and ideally label free screening
50
51 technologies. It is crucial, however, that these models and technologies are well validated against
52
53 established gold standards ref^{33,34}.
54
55
56
57
58
59
60

1
2
3 Raman spectra, in principle, contain high content information about the biochemical make
4 up of the sample, and changes to it, related to pathology or an external agent. Raman spectra
5 contain numerous peaks which vary dependently and independently of each other. Crucially, for real
6 applications and particularly in the instance of drug interactions, it is difficult to tell whether these
7 differences are inherently based on cell to cell variability or whether they are dependent on the
8 primary action of the drug (i.e. the direct chemical effects) or the secondary effects the drug has on
9 the cell (i.e. the response of the cell to said drug).
10
11
12
13
14
15
16
17
18

19 In this study, simulated datasets were used to evaluate the capability of PLSR to extract
20 known and systematic spectral variation from a control dataset, which contained intrinsic
21 experimental variability. The spectral variations introduced varied linearly with the applied drug
22 dose and also with the measured cell population response, as measured by a standard cytotoxicity
23 assay. Notably, however, the two spectral variations are not completely independent, as the viability
24 response is sigmoidally dependent on the applied dose.
25
26
27
28
29
30
31
32

33 In the case where only a concentration dependent systematic variation in the spectra is
34 introduced, the PLSR model provides an accurate predictive response tool, the regression co-
35 efficients of which are based on the systematic variation which has been introduced to the dataset,
36 linearly dependent on the targets. The model shows high sensitivity, and the limits of detection are
37 determined only by the intrinsic variability of the experimental method, as determined by the PLSR
38 of the Control spectral dataset. This limit can be improved by optimising sample preparation and
39 measurement protocols. In principle, such a PLSR model can predict the response of a drug dose in a
40 cell population, or determine an unknown drug dose from a measured spectral response.
41
42
43
44
45
46
47
48
49

50 However, the spectral changes which result from the interaction and action of a drug within
51 a cell are manifold, and it is of interest to differentiate the spectral signatures of the direct
52 interaction from the subsequent cellular response. Notably, this study demonstrates that, although
53 PLSR predictive models based on regression of the combined dataset, including all spectral
54
55
56
57
58
59
60

1
2
3 responses, against the target of concentration range produce a similarly accurate, linear predictive
4
5 model, the contributing RCs are only derived exclusively from the introduced concentration
6
7 dependent variations in ranges where all other spectral variations are limited. For example, as
8
9 shown in Figures 8 and 9, regression over the limited range of C+4 produces a model which is based
10
11 on RCs which includes contributions derived from the direct effect of the interaction of the drug
12
13 within the cell (Concentration construct), as well as the resultant cytological response (Viability
14
15 construct). Thus, care should be taken in interpreting the spectral features which contribute to such
16
17 regressions to elucidate the underlying mechanisms.
18
19

20
21 Nevertheless, in sublethal regions, the direct effects of the drug interaction can confidently
22
23 be investigated employing such a PLSR analysis of Raman spectral data, independent of the
24
25 cytological responses, and these are easily discernible above the intrinsic variability of the control.
26
27 Although this seems a trivial conclusion, such rapid, label free analysis could prove invaluable in
28
29 screening of, for example, the mechanisms and efficacy of drug interactions, evaluating drug uptake
30
31 and receptor binding²⁵ or nanoparticle uptake and trafficking in regions where cytotoxicity assays are
32
33 insensitive.
34
35

36
37 The use of a parallel cytotoxic assay such as MTT serves as a range finding test to establish
38
39 the IC₅₀, but also provides vital information about the sublethal doses and maximum responses. It
40
41 also provides a target for regression of the data in the regions of toxicity. Thus, the subsequent
42
43 cytological effects can be differentiated from the direct chemical effects of the agent and extracted
44
45 from the overall spectral response in the dose range where the viability is impacted, and the cellular
46
47 response can be independently mapped spectroscopically, as a function of dose and time. Notably,
48
49 the model described here, which includes a single spectral construct to represent the cellular
50
51 response is very simplistic, as the response is a cascade of many responses, depending on the
52
53 mechanism of interaction³⁵. Alternative cytological gold standard assays for cancer, such as the
54
55 sulphorhodamine B assay, and human tumour cell lines such as NCI60 human tumour cell line,
56
57
58
59
60

1
2
3 should also be considered to broaden the model^{36,37}. Nevertheless, the analysis presented here
4
5 demonstrates that the spectral fingerprints of the direct mechanisms of interaction and the
6
7 subsequent cellular responses can be independently extracted from the dose dependent spectral
8
9 data, and thus, ultimately with improved screening sensitivities and speeds, Raman spectroscopy
10
11 could be employed to monitor in quasi realtime, in a lable free manner, the efficacy and mode of
12
13 action of, for example chemotherapeutic agents and other exogenous agents, laying the basis for
14
15 improved quantitative structure activity relationships to guide drug development or chemical
16
17 regulation strategies.
18
19

20 21 **Conclusions**

22
23 This study demonstrates the reliability and also limitations of PLSR as a method for predictive
24
25 modelling and analysis of spectroscopic signatures of cellular responses to exogeneous agents such
26
27 as radiation, chemotherapeutic agents or toxins. The spectroscopic profiles at any dose/time point
28
29 can derive from a complex mixture of direct interactions within the cell and a cascade of subsequent
30
31 cellular response. The analysis demonstrates that care should be taken in choosing the response
32
33 range and also highlights the importance of parallel cytological assays in guiding the modelling and
34
35 analysis. Correct choice of range can help differentiate between the signatures of direct interactions,
36
37 which are dominant at sub-lethal doses and those of the subsequent cellular response which evolve
38
39 with increasing dose.
40
41
42
43

44 The study also demonstrates the importance of simulated datasets in exploring the potential
45
46 as well as the limits of the analytical techniques. Notably, the use of real experimental data which
47
48 contains sample variability and instrumental response factors as a basis of the simulated dataset
49
50 helps to visualise the lower limits of sensitivity.
51
52

53 The results indicate that Raman spectroscopic screening combined with such regression
54
55 models and feature selection techniques, in parallel with conventional cytotoxicity assays, can be
56
57
58
59
60

1
2
3 used to screen for the efficacy of drug interactions and can contribute to understanding the
4
5 mechanisms of interaction.
6
7
8
9

10
11 **Acknowledgement:** This research was supported by the Integrated NanoScience Platform, Ireland
12
13 (INSPIRE), funded under the Higher Education Authority PRTL (Programme for Research in Third
14
15 Level Institutions) Cycle 5, co-funded by the Irish Government and the European Union Structural
16
17 fund, and Science Foundation Ireland (08/PI/11).
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

1. F. M. Lyng, E. O. Faoláin, J. Conroy, a D. Meade, P. Knief, B. Duffy, M. B. Hunter, J. M. Byrne, P. Kelehan, and H. J. Byrne, *Exp. Mol. Pathol.*, 2007, **82**, 121–9.
2. I. Taleb, G. Thiéfin, C. Gobinet, V. Untereiner, B. Bernard-Chabert, A. Heurgué, C. Truntzer, P. Hillon, M. Manfait, P. Ducoroy, and G. D. Sockalingum, *Analyst*, 2013, **138**, 4006–14.
3. P. Crow, B. Barrass, C. Kendall, M. Hart-Prieto, M. Wright, R. Persad, and N. Stone, *Br. J. Cancer*, 2005, **92**, 2166–70.
4. T. J. Harvey, E. Gazi, A. Henderson, R. D. Snook, N. W. Clarke, M. Brown, and P. Gardner, *Analyst*, 2009, **134**, 1083–91.
5. F. Bonnier, P. Knief, B. Lim, a D. Meade, J. Dorney, K. Bhattacharya, F. M. Lyng, and H. J. Byrne, *Analyst*, 2010, **135**, 3169–77.
6. K. Klein, A. M. Gigler, T. Aschenbrenner, R. Monetti, W. Bunk, F. Jamitzky, G. Morfill, R. W. Stark, and J. Schlegel, *Biophys. J.*, 2012, **102**, 360–8.
7. M. Miljković, T. Chernenko, M. J. Romeo, B. Bird, C. Matthäus, and M. Diem, *Analyst*, 2010, **135**, 2002–13.
8. C. Matthäus, T. Chernenko, J. a Newmark, C. M. Warner, and M. Diem, *Biophys. J.*, 2007, **93**, 668–73.
9. P. Bassan, A. Sachdeva, A. Kohler, C. Hughes, A. Henderson, J. Boyle, J. H. Shanks, M. Brown, N. W. Clarke, and P. Gardner, *Analyst*, 2012, **137**, 1370–7.
10. H. Nawaz, F. Bonnier, P. Knief, O. Howe, F. M. Lyng, A. D. Meade, and H. J. Byrne, *Analyst*, 2010, **135**, 3070–6.
11. H. Nawaz, F. Bonnier, A. D. Meade, F. M. Lyng, and H. J. Byrne, *Analyst*, 2011, **136**, 2450–63.
12. T. Chernenko, R. R. Sawant, M. Miljkovic, L. Quintero, M. Diem, and V. Torchilin, *Mol. Pharm.*, 2012, **9**, 930–6.
13. J. Dorney, F. Bonnier, A. Garcia, A. Casey, G. Chambers, and H. J. Byrne, *Analyst*, 2012, **137**, 1111–9.
14. M. E. Keating, F. Bonnier, and H. J. Byrne, *Analyst*, 2012, **137**, 5792–802.
15. P. Bassan, A. Kohler, H. Martens, J. Lee, H. J. Byrne, P. Dumas, E. Gazi, M. Brown, N. Clarke, and P. Gardner, *Analyst*, 2010, **135**, 268–77.
16. H. Byrne, K. Ostrowska, and H. Nawaz, *Opt. Spectrosc. Comput. Methods Biol. Med.*, 2014, **14**, 355–399.
17. M. Miljković, T. Chernenko, M. J. Romeo, B. Bird, C. Matthäus, and M. Diem, *Analyst*, 2010, **135**, 2002–13.

- 1
- 2
- 3 18. A. D. Meade, H. J. Byrne, and F. M. Lyng, *Mutat. Res.*, 2010, **704**, 108–14.
- 4
- 5 19. K. M. Ostrowska, A. Malkin, A. Meade, J. O’Leary, C. Martin, C. Spillane, H. J. Byrne, and F. M.
- 6 7 Lyng, *Analyst*, 2010, **135**, 3087–93.
- 8
- 9 20. R. M. Balabin and S. V Smirnov, *Anal. Chim. Acta*, 2011, **692**, 63–72.
- 10
- 11 21. M. Jimenez-Hernandez, C. Hughes, P. Bassan, F. Ball, M. D. Brown, N. W. Clarke, and P.
- 12 13 Gardner, *Analyst*, 2013, **138**, 3957–66.
- 14 22. K. W. C. Poon, F. M. Lyng, P. Knief, O. Howe, A. D. Meade, J. F. Curtin, H. J. Byrne, and J.
- 15 16 Vaughan, *Analyst*, 2012, **137**, 1807–14.
- 17 23. H. Nawaz, A. Garcia, A. D. Meade, F. M. Lyng, and H. J. Byrne, *Analyst*, 2013, **138**, 6177–84.
- 18
- 19 24. D. Rohleder, W. Kiefer, and W. Petrich, *Analyst*, 2004, **129**, 906–11.
- 20
- 21 25. J. Black and P. Leff, *Proc R Soc L. B Biol Sci.*, 1983, **220**, 141–162.
- 22
- 23 26. S. Wold, M. Sjöström, and L. Eriksson, *Chemom. Intell. Lab. ...*, 2001, 109–130.
- 24
- 25 27. A. Meade, C. Clarke, H. Byrne, and F. Lyng, *Radiat. Res.*, 2010, **2**, 225–37.
- 26
- 27 28. K. Vermuza and P. Flizmoser, *Introduction to Multivariate Statistical Analysis in*
- 28 29 *Chemometrics*, CRC Press, 2009.
- 30
- 31 29. H. Martens and T. Næs, *Multivariate Calibration*, John Wiley & Sons, 1994.
- 32
- 33 30. THE EUROPEAN PARLIAMENT AND THE COUNCIL OF THE EUROPEAN UNION, *Off. J. Eur.*
- 34 35 *Union*, 2010, 33–79.
- 36
- 37 31. U. S. Congress, 2001, 2721–2725.
- 38
- 39 32. W. Russell, R. Burch, and C. Hume, *The principles of humane experimental technique*,
- 40 41 Methuen, London, 1959.
- 42
- 43 33. A. Tfayli, F. Bonnier, Z. Farhane, D. Libong, H. J. Byrne, and A. Baillet-Guffroy, *Exp. Dermatol.*,
- 44 45 2014, **23**, 441–3.
- 46 34. F. Bonnier, M. Keating, T. Wróbel, K. Majzner, M. Baranska, A. Garcia, A. Blanco, and H. J.
- 47 48 Byrne, *Toxicol. Vit.*, 2014, **29**, 124–131.
- 49 35. M. Maher, P. C. Naha, S. P. Mukherjee, and H. J. Byrne, *Toxicol. Vit.*, 2014, **28**, 1449–60.
- 50
- 51 36. V. Vichai and K. Kirtikara, *Nat. Protoc.*, 2006, **1**, 1112–1116.
- 52
- 53 37. R. Shoemaker, *Nat. Rev. Cancer*, 2006, **6**, 813–823.
- 54
- 55
- 56
- 57
- 58
- 59
- 60