



**QUEEN'S
UNIVERSITY
BELFAST**

Testing for multivariate normality in mass spectrometry imaging data: a robust statistical approach for clustering evaluation and the generation of synthetic mass spectrometry imaging datasets

Dexter, A., Race, A. M., Styles, I., & Bunch, J. (2016). Testing for multivariate normality in mass spectrometry imaging data: a robust statistical approach for clustering evaluation and the generation of synthetic mass spectrometry imaging datasets. *Analytical Chemistry*, 88(22), 10893-10899.
<https://doi.org/10.1021/acs.analchem.6b02139>

Published in:
Analytical Chemistry

Document Version:
Peer reviewed version

Queen's University Belfast - Research Portal:
[Link to publication record in Queen's University Belfast Research Portal](#)

Publisher rights

Copyright 2016 American Chemical Society.
This work is made available online in accordance with the publisher's policies. Please refer to any applicable terms of use of the publisher.

General rights

Copyright for the publications made accessible via the Queen's University Belfast Research Portal is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The Research Portal is Queen's institutional repository that provides access to Queen's research output. Every effort has been made to ensure that content in the Research Portal does not infringe any person's rights, or applicable UK laws. If you discover content in the Research Portal that you believe breaches copyright or violates any law, please contact openaccess@qub.ac.uk.

Open Access

This research has been made openly available by Queen's academics and its Open Research team. We would love to hear how access to this research benefits you. – Share your feedback with us: <http://go.qub.ac.uk/oa-feedback>

Testing for multivariate normality in mass spectrometry imaging data: A robust statistical approach for clustering evaluation and the generation of synthetic mass spectrometry imaging datasets

Alex Dexter^{1,2}, Alan M. Race², Iain B. Styles³, Josephine Bunch^{2,4}.

¹PSIBS Doctoral Training Centre, University of Birmingham, Edgbaston, Birmingham, B15 2TT, United Kingdom; ²National Physical Laboratory, Teddington, Middlesex TW11 0LW, UK; ³

School of Computer Science, University of Birmingham, Edgbaston, Birmingham, B15 2TT, United Kingdom; ⁴School of Pharmacy, University of Nottingham, Nottingham, Nottinghamshire NG7

2RD, UK

Abstract

Spatial clustering is a powerful tool in mass spectrometry imaging (MSI), and has been demonstrated to be capable of differentiating tumour types, visualising intra-tumour heterogeneity, and segmenting anatomical structures. Several clustering methods have been applied to mass spectrometry imaging data but a principled comparison and evaluation of different clustering techniques presents a significant challenge. We propose that testing whether the data has a multivariate normal distribution within clusters can be used to evaluate the performance when using algorithms that assume normality in the data such as k-means clustering. In cases where clustering has been performed using the cosine distance, conversion of the data to polar coordinates prior to normality testing should be performed to ensure normality is tested in the correct coordinate system. In addition to these evaluations of internal consistency, we demonstrate that the multivariate normal distribution can then be used as a basis for statistical modelling of MSI data. This allows the generation of synthetic MSI datasets with known ground truth, providing a means of external clustering evaluation. To demonstrate this, reference data from seven anatomical regions of an MSI image of a coronal section of mouse brain were modelled. From this a set of synthetic data based on this model was generated. Results of r^2 fitting of the chi-squared quantile-quantile plots on the seven anatomical regions confirmed that the data acquired from each spatial region was found to be closer to normally distributed in polar space than in Euclidean. Finally, principal component analysis was applied to a single dataset which included synthetic and real data. No significant differences were found between the two data types indicating the suitability of these methods for generating realistic synthetic data.

Introduction

Data mining is a valuable tool in mass spectrometry imaging (MSI), where even a single image can contain more information than can be feasibly interpreted by a single person in a realistic timeframe. Often, a few m/z values or pixels of interest are selected for analysis based on known information about the sample. It is becoming increasingly clear however that simple univariate analysis is both impractical and does not take full advantage of the rich content of the data, and that multivariate analysis methods are increasingly important to effectively mine this data.¹⁻³ One of the main tasks for which multivariate analysis is used in MSI is to segment different regions of an image for the purpose of diagnosis of diseases or to improve disease understanding, and to segment anatomical regions for comparison to histology in order to more fully understand the molecular composition of different anatomical regions.^{4,5}

Clustering techniques divide the dataset into classes and assign a single class label to each pixel and as such provide a clear categorisation of the data. However, the idea of a cluster of data is arbitrary, relying on the notion of “similarity” which can be formulated in many ways. There are many clustering techniques, each of which makes specific assumptions about the data, and will therefore categorise a given dataset very differently depending on the validity of the assumption.⁶⁻¹¹ There is no *a priori* method for determining which method is appropriate for a given dataset. A further and very significant challenge to clustering in MSI is the size of the data itself, both in terms of the number of data points and the dimensionality. A number of different clustering algorithms have been applied to MSI data^{1-3,5,12}, each of which makes specific assumptions about the properties of the data, and has inherent advantages and disadvantages.^{1-3,5,12} Due to its simplicity, relatively low computational requirements⁷, and wide availability in many different languages², k -means clustering is one of the most popular algorithms for clustering in MSI.^{2,13,14} This can distinguish between anatomies within mouse brain tissue¹³, distinguish tumour margins¹⁵ and even intra tumour heterogeneity⁴. Given a set of spectra, k -means clustering aims to partition the n spectra into k sets so as to minimize the intra-cluster sum of distances of each point in the cluster to its cluster centre. An illustration of the iterative process of k -means clustering is provided in Figure S1 in the Supporting Information.

The distance metric used by the clustering algorithms to compare one spectrum to another (Figure 1), and any normalisation strategies applied to the data prior to analysis have a significant effect on the results. In MSI, there can be significant variations in the data that are derived from a number of different experimental sources. For example, variability in sample preparation¹⁶, and laser instability¹⁷ both introduce a source of non-biological variance within the data. Minimising these effects by normalisation is common but does not and cannot remove all non-biological variations.¹⁸ Nevertheless, normalisation of the data, or pseudo normalisation achieved by the use of the cosine distance, reduce the effects of these variations, and thereby improve the clustering results. In the commonly applied TIC normalisation, each spectrum is normalised to have unit sum intensity (also referred to as l_1 norm). The cosine similarity is also intensity-independent and therefore also has potential to reduce the impact of some of these variations on clustering performance (Figure 1b).

Most applications of k -means clustering in MSI have used the Euclidean distance metric^{2,13,14}, and where normalisation has been used, total ion count (TIC) normalisation is most common.^{4,19} Most attempts to evaluate clustering results in MSI have used manual examination or comparison to complimentary modalities such as histological analysis.⁵ Recently Oetjen *et al.* published a series of benchmark 3D datasets with histological information²⁰; however the limited chemical information provided by histology means that segmentations do not always match chemical information provided by MSI.⁴

There are many different methods for quantitatively evaluating the success of clustering, which can be divided into two types; internal and external. Internal evaluation uses the intrinsic properties of the clustering result, usually by comparing the data within each cluster to the data outside of the cluster.²¹ Previous attempts to evaluate clustering in MSI have used internal evaluation measures but these have proven inconclusive at best.^{22,23} External evaluation on the other hand compares the clustering results to known ground truths such that true and false positives and negative can be computed. Using this information, values such as sensitivity and specificity can be calculated

alongside validation measures such as the Rand and Jaccard indices.²⁴ Since the comparison is to known information there is no concern of bias towards a given algorithm or distance metric and so can be used as a method for accurately and reliably comparing and evaluating clustering algorithms or workflows. The main limitation of external evaluation is the need for a ground truth to compare against. Since MSI is generally used as an exploratory tool, usually on biological samples, most datasets will not have a ground truth and thus these external evaluations are usually not possible.²⁵

One of the primary assumptions of the *k*-means clustering and other algorithms is that the data within clusters is normally distributed. Previously, in other fields, methods have been used to evaluate whether the data within clusters is normally distributed to evaluate the clustering performance²⁶, or to determine whether to continue to divide clusters further^{27,28}. By evaluating the degree of normality within the clusters, when clustering with an algorithm that assumes normality, it is possible to evaluate how well the data fits this assumption and thus how appropriate it is. For univariate data, normality testing is relatively straightforward, and there are a number of tests for normality such as Shapiro-Wilks²⁹, Kolmogorov-Smirnov³⁰, and Cramer-Von Mises³¹ tests. This is more challenging in multivariate data since there will be many dimensions each with different variance and means.³² It is possible to test for multivariate normality however using quantile-quantile plots.³³ If the data is multivariate normal then the Mahalanobis distance will have a χ_p^2 distribution.³⁴ Therefore plotting the Mahalanobis distance from each pixel to its relevant distribution versus a χ_p^2 distribution where *p* is the dimensions of the data will give a straight line if the data is multivariate normal.

In this work we show how multivariate normality testing can be used to evaluate the appropriateness of difference distance metrics in *k*-means clustering. We also show how the multivariate normal model can be used as a basis for generating synthetic mass spectrometry imaging datasets, thereby providing samples with a ground truth against which to quantitatively evaluate multivariate analysis methods in MSI, as well as other computational analysis methods.

Materials and methods

Image acquisition; Coronal mouse brain was sectioned to 12 μm thickness and thaw mounted onto glass slides (Superfrost, Thermo Fisher Scientific, Waltham, MA USA), before being coated with α -cyano-4-hydroxycinnamic acid (CHCA) matrix (5 mg/mL, 80% MeOH 0.1% TFA) using an automated pneumatic sprayer (TM-sprayer, HTX imaging, Chapel Hill, NC, USA). Matrix-assisted laser desorption/ionisation (MALDI) images were acquired using a Synapt G2Si (Waters, Manchester, UK), using a pixel size of 45 μm in both x and y , and an m/z range of 100-1200 Da.

Data processing and analysis; Data processing was performed on an Intel Xeon quad core CPU E5-2637 v2 (3.50 GHz) with 64 GB of RAM. All data were converted from proprietary format to the mzML format using msconvert as part of ProteoWizard³⁵ software then into imzML using imzMLConverter³⁶. This was then imported into MATLAB (version R2014a and statistics toolbox, The Math-Works, Inc., Natick, MA, USA) using Spectral Analysis, an in-house mass spectrometry imaging software. k -means clustering was performed using the function *kmeans* from the Matlab Statistics toolbox using the parameters specified in the upcoming experiments and three replicates and random starting clusters. Normality testing was performed on the data within each cluster by plotting the squared Mahalanobis distance from each pixel to the distribution within its cluster against a chi-square distribution with a number of degrees of freedom equal to the dimensions of the data.³³ The Mahalanobis distance for the data within each cluster was calculated by first performing PCA, removing components with zero variance, and scaling such that each component has a standard deviation of 1; then calculating the squared Euclidean distance of each pixel to the mean of its assigned cluster.³⁷ For data clustered using the cosine distance metric, data were first converted into a polar coordinate system, comprising of a distance from the origin r , and a series of angles from the origin θ_{n-1} relative to each of the coordinate axes where n is the dimensionality of the data.³⁸ The angles from the co-ordinate axes were then used to determine normality of the angular distribution. For creating the plots the Mahalanobis distance and chi-squared values were all rescaled to between 0 and 1 in order to plot them all on a common axis.

Synthetic data were generated using the following workflow;

1. Convert the reference data to polar coordinates
2. Test for normality of the reference data in polar coordinates using chi-squared quantile-quantile plotting/
3. If the reference data is multivariate normal then calculate the means and covariances of the reference data in polar coordinates.
4. Generate a set of synthetic multivariate normally distributed data with the mean and covariance of the reference data using the *mvrnd* function in MATLAB
5. Convert the synthetic multivariate normal data back to Cartesian coordinates
6. Populate a spatial mask with the synthetic data

Results and Discussion

Two synthetic datasets were generated to simulate data that is normally distributed in Euclidean and polar coordinates respectively. Clustering was performed using *k*-means with both the Euclidean and cosine distances. Normality testing via quantile-quantile plots revealed that both synthetic datasets clustered well under the cosine distance, whereas the data distributed normally in polar coordinates did not cluster well when the Euclidean distance was used (Figure 2).

k-means clustering was then performed on an MSI image from coronal mouse brain with $k = 2-10$, 7 clusters were then chosen based on visual assessment of the resulting images, along with comparison to the Allen brain atlas. When applied to an MSI image of a biological system (coronal mouse brain), the Chi squared quantile plots show that the data within clusters obtained using the cosine distance have a higher r^2 value than the data within the clusters using the Euclidean distance (average r^2 0.99 compared to 0.87 Figure 3 A and B). This means that the data in the clusters formed using the cosine distance are closer to normally distributed than the Euclidean distance. This indicates that the cosine distance is the more appropriate distance metric for cluster with on this dataset based on the multivariate normal assumption of the *k*-means algorithm. The inappropriateness of *k*-means with the Euclidean distance in this case mirrors the visually poor results obtained with respect to the anatomical features expected from coronal mouse brain as seen in the Allen brain atlas (Figure 3).³⁹ In

comparison the cosine distance gives visually clearer results, and the distribution of points within clusters are more normally distributed in the appropriate space. We note that use of the common TIC normalisation *decreases* the normality of the data, and does not produce visually clearer segmentation images (Figure 3c). The reason for this is that TIC normalisation rescales all data points such that they lie on the surface of a hyperdiamond (lines of constant L_1 norm). Thus, they are certainly not normally distributed as one dimension is condensed. They might be normally distributed if you consider only positions on the hypersurface. Results obtained from additional datasets (sagittal rat brain and mouse lung tissue) produce similar results with respect to comparison of normality to visual appearance of clustering results and are provided in the supplementary material (Figures S2 and S3). It is worth noting that the values produced from the r^2 fitting cannot easily be directly interpreted, as it will be dependent on the number of data points, and the dimensionality of the data. Therefore it is recommended as a means to compare results, and caution should be taken when inferring additional information from them.

It is also worth noting that the shape of the q-q plots are not completely linear, a feature that can arise from a number of different sources. For example, the presence of a few outliers will skew the distribution towards a sigmoidal shape as is observed when the cosine distance is used (Figures S4 and S5). This is caused by the outliers skewing the mean of the data, and thus altering the Mahalanobis distance for every point. While this effect is minimised through the variance scaling process, some effect can still be observed. Alternately, a circular distribution of data with a core of normal data within produces a similar shaped, apparent bilinear plot to those observed when using the Euclidean distance (Figure S6). This is not indicative of two normally distributed sets of data however which produces a different shaped plot (Figure S8). For further examples of how other distributions of data will affect these plots see figures S4 to S12. However, we note that caution is required when generalising from these plots from two dimensional data into the higher dimensional space in which MSI data sits.

While distance metric determination is a crucial factor in any clustering algorithm, there are still many other parameters which must also be selected such as the number of clusters, and the

method for centroid initiation. In addition to this, there are many other clustering approaches such as density based clustering which do not assume multivariate normality in the data. Therefore, a method to generate datasets with a ground truth is required to assess the suitability of these approaches and to permit a comparison of different clustering approaches. Data simulated from first principles is one approach that is used to achieve this in other fields. However, while some aspects of image formation and noise in MSI are well understood there are still a large number of unknowns in aspects such as sample preparation and ionisation.^{40,41} One approach is to take existing peak lists and to then simulate the known variables and apply these to this peak list.⁴² A robust method is needed however to generate peak lists that are well controlled, but still representative of the variance expected from biological samples. A new biological sample could be analysed each time a new set of spectra are required, but using new animal or human tissue each time a different number of regions or pixels is required is neither practical or ethical. In other areas such as financial prediction, and geological analysis, statistical modelling is used to convert discrete data into a continuous function, thereby allowing resampling to generate the desired number of data points. Statistical modelling assumes that data from a population are derived from a known probability distribution function. Provided that the model adequately describes the data, the underlying distribution can then be resampled to give a new synthetic dataset with any desired number of data points. This new synthetic dataset will have the same distribution as the original reference dataset that the model was derived from. For large and high dimensional data, model generation and parameter estimation can be challenging, however, the multivariate normal model parameters can be easily estimated even for very large data.⁴³ As previously demonstrated, clustered MSI data closely approximates to a multivariate normal distribution when the data is converted to polar coordinates. This means that the multivariate normal distribution can be used as the basis for statistical modelling for MSI data. The small deviations from normal are most probably due to a few outlier pixels within the data, rather than a deviation from normal within the majority of the data itself. This is demonstrated by the similarity of the sigmoidal nature of the plot (Figure 3B) to that generated using simulated outliers (Figures S4 and S5). This suggests that synthetic spectra generated in this way are representative of those observed in the real

data, and thereby serve as a basis to introduce and explore additional experimental or instrumental variabilities in a controlled way.

In order to perform statistical modelling of MSI data, a series of seven anatomical features from an MSI image of the previously shown mouse brain were used as a reference dataset (Figure 4). These regions were generated based on the analysis of selected ion images and PCA scores, in comparison to a high resolution optical image (Figure S13), and the Allen brain atlas.³⁹ These datasets were then tested for normality in polar and Euclidean space using the chi-squared quantile plots shown previously and showed a high degree of normality throughout polar space but not Euclidean (Table 1). This suggests that the multivariate normal model can be used to summarise the properties of this data. From this model, a new synthetic dataset was generated by resampling from the distribution with the same number of pixels as the original reference data. The synthetic spectra from a number of the different anatomical regions were then visually compared to the original reference spectra (Figure 5). The synthetic and real spectra show a high degree of spectral similarity, and expected features such as isotope ratios and fragments are preserved, thus ensuring the realism of the synthetic data. Some differences in the spectra are observed, since the synthetic spectra are sampled from a distribution and will therefore contain the same underlying variance as the reference data. This is important since biological samples vary, and so in order to be realistic, the synthetic data must incorporate this variance.

Visual comparison of the spectra is insufficient to evaluate datasets that will be analysed by multivariate methods. Therefore, in order to evaluate how closely the synthetic data matches real data, a new dataset, comprising of both synthetic and real spectra was generated. Principal component analysis was then performed on this combined dataset to determine if the statistical modelling process introduced any additional observable variance. No principal component scores were found to separate the synthetic from real data (Figure 6 and S14). This means that even when all mass channels are considered, the difference between the synthetic and real data is smaller than that between different anatomies or the spectral noise within the data and supports the suggestion that the differences from normal are likely to be outlier pixels. As such the statistical modelling of appropriately segmented

MSI data using a multivariate normal distribution can generate realistic spectra in order to create new datasets with known ground truth for external evaluation of clustering in mass spectrometry imaging.

Large synthetic datasets can also be generated rapidly using this approach, by simply taking more samples from the multivariate normal distribution. To demonstrate this, a dataset containing nine times the number of pixels of the original reference data was generated (187,452 pixels from 20,825 in the original). This represents this size of data from an area three times the size in each dimension, or if the image had been acquired with 15 μm rather than 45 μm pixels. These new data were generated in approximately 5 minutes, but it would have required around 36 hours to acquire the same number of pixels experimentally. PCA performed on a combined dataset containing the new larger dataset and the original reference data still shows no separation between the synthetic and real data, demonstrating that this approach scales to large datasets without any statistically detectable changes occurring in the data (Figure S15). While in both these cases the full seven regions were used to generate synthetic data, an image containing any desired number of regions can be generated using this approach, provided there is a suitable set of reference data. This means that the performance of different clustering algorithms or multivariate analysis methods can be evaluated with respect to the size and complexity of the data in terms of expected features. In addition, no new tissue sections are required, allowing the potential to minimise animal usage in computational studies in MSI. We note that the synthetic images appear more speckled than the reference data. This is because when populating the spatial masks with spectra, no spatial smoothing is applied and neighbouring pixels are statistically independent. This could potentially be overcome by also maximising the similarity of neighbouring pixel, but for clustering evaluation this is unnecessary.

Conclusions

Robust evaluation of clustering in MSI allows us to understand its limitations and what can be deduced from its results. In the case of k -means clustering and other algorithms that assume normality of the data (such as agglomerative hierarchical clustering with Ward's linkage), we have shown that in the absence of ground-truth data, evaluation of multivariate normality in the intra cluster

distributions is an internal test that can be used on MSI data to determine, post-analysis, whether clustering should be performed using the cosine or Euclidean distance. Where possible, external evaluation methods should be used when comparing novel algorithms or parameters, using a ground truth that is representative of samples of interest. We have demonstrated that synthetic data generated by statistical modelling is a suitable means to achieve this. In addition this approach allows large datasets to be generated rapidly allowing evaluation and comparison of both existing and new methods as the data increases in size.

Acknowledgements

This research was funded by NPL strategic research programmes 116301 and 117194. The authors thank Dr Rory Steven for his assistance with MALDI MS imaging experiments, and Professor Helen Cooper for useful discussions. AD gratefully acknowledges financial support from the EPSRC through a studentship from the PSIBS Doctoral Training Centre (EP/F50053X/1) in collaboration with the National Physical Laboratory and AstraZeneca. Data supporting this research is openly available from the University of Birmingham data archive at <http://findit.bham.ac.uk/>.

References

- (1) Muir, E.; Ndiour, I.; Le Goasduff, N.; Moffitt, R. A.; Liu, Y.; Sullards, M. C.; Merrill, A. H.; Chen, Y.; Wang, M. D. *In Bioinformatics and Bioengineering*, Proceedings of the 7th IEEE International Conference on; IEEE, 2007, pp 472-479.
- (2) McCombie, G.; Staab, D.; Stoeckli, M.; Knochenmuss, R. *Anal.Chem.* **2005**, *77*, 6118-6124.
- (3) Fonville, J. M.; Carter, C. L.; Pizarro, L.; Steven, R. T.; Palmer, A. D.; Griffiths, R. L.; Lalor, P. F.; Lindon, J. C.; Nicholson, J. K.; Holmes, E. *Anal. Chem.* **2013**, *85*, 1415-1423.
- (4) Jones, E. A.; van Remoortere, A.; van Zeijl, R. J.; Hogendoorn, P. C.; Bovée, J. V.; Deelder, A. M.; McDonnell, L. A. *PLoS One* **2011**, *6*, e24913.
- (5) Deininger, S. r.-O.; Ebert, M. P.; Fütterer, A.; Gerhard, M.; Röcken, C. *J. Proteome Res.***2008**, *7*, 5230-5236.
- (6) Estivill-Castro, V. *SIGKDD Explor.* **2002**, *4*, 65-75.
- (7) Hartigan, J. A.; Wong, M. A. *Appl. Stat.* **1979**, 100-108.
- (8) Jain, A. K. *Pattern Recogn. Lett.* **2010**, *31*, 651-666.
- (9) Birant, D.; Kut, A. *Data Knowl. Eng.* **2007**, *60*, 208-221.
- (10) Fu, L.; Medico, E. *BMC Bioinf.* **2007**, *8*, 3.
- (11) Choong, M. Y.; Kow, W. Y.; Chin, Y. K.; Angeline, L.; Teo, K. T. K. *Control System, Computing and Engineering*, Proceedings of the IEEE International Conference on, 2012, p 430-435.
- (12) Trede, D.; Schiffler, S.; Becker, M.; Wirtz, S.; Steinhorst, K.; Strehlow, J.; Aichler, M.; Kobarg, J. H.; Oetjen, J.; Dyatlov, A. *Anal.Chem.* **2012**, *84*, 6079-6087.
- (13) Race, A. M.; Steven, R. T.; Palmer, A. D.; Styles, I. B.; Bunch, J. *Anal. Chem.* **2013**, *85*, 3071-3078.
- (14) Alexandrov, T.; Becker, M.; Deininger, S. O.; Ernst, G.; Wehder, L.; Grasmair, M.; von Eggeling, F.; Thiele, H.; Maass, P. *J. Proteome Res.***2010**, *9*, 6535-6546.
- (15) Alexandrov, T.; Becker, M.; Guntinas-Lichius, O.; Ernst, G.; von Eggeling, F. *J. Cancer Res. Clin. Oncol.* **2013**, *139*, 85-95.
- (16) Goodwin, R. J. *Proteomics* **2012**, *75*, 4893-4911.
- (17) Steven, R. T.; Dexter, A.; Bunch, J. *Methods* **2016**.
- (18) Deininger, S.-O.; Cornett, D. S.; Paape, R.; Becker, M.; Pineau, C.; Rauser, S.; Walch, A.; Wolski, E. *Anal. Bioanal. Chem.* **2011**, *401*, 167-181.
- (19) Abdelmoula, W. M.; Carreira, R. J.; Shyti, R.; Balluff, B.; van Zeijl, R. J.; Tolner, E. A.; Lelieveldt, B. F.; van den Maagdenberg, A. M.; McDonnell, L. A.; Dijkstra, J. *Anal.Chem.* **2014**, *86*, 3947-3954.
- (20) Oetjen, J.; Veselkov, K.; Watrous, J.; McKenzie, J. S.; Becker, M.; Hauberg-Lotte, L.; Kobarg, J. H.; Strittmatter, N.; Mróz, A. K.; Hoffmann, F. *GigaScience* **2015**, *4*, 1-8.
- (21) Liu, Y.; Li, Z.; Xiong, H.; Gao, X.; Wu, J. *In Data Mining*, Proceedings of the IEEE 10th International Conference on; IEEE, 2010, pp 911-916.
- (22) Van de Plasa, R.; Ojedaa, F.; Dewile, M.; Van, L.; Den Bosche, B. D. M.; Waelkensbcd, E. *Bioinformatics* **2006**.
- (23) Sarkari, S.; Kaddi, C. D.; Bennett, R. V.; Fernandez, F. M.; Wang, M. D. *In Engineering in Medicine and Biology Society*, Proceedings of the 36th Annual International Conference of the IEEE **2014**, pp 4771-4774.
- (24) Rand, W. M. *J. Am. Stat. Assoc.* **1971**, *66*, 846-850.
- (25) Garden, R. W.; Sweedler, J. V. *Anal. Chem.* **2000**, *72*, 30-36.
- (26) Mao, J.; Jain, A. K. *Neural Networks*, IEEE Transactions on **1996**, *7*, 16-29.
- (27) Hamerly, G.; Elkan, C. *Adv. Neural Inf. Process. Syst.* **2004**, *16*, 281.
- (28) Steinley, D. *Br. J. Math. Stat. Psychol.* **2006**, *59*, 1-34.
- (29) Shapiro, S. S.; Wilk, M. B. *Biometrika* **1965**, 591-611.
- (30) Lilliefors, H. W. *J. Am. Stat. Assoc.* **1967**, *62*, 399-402.
- (31) Darling, D. A. *Ann. Math. Stat.* **1957**, 823-838.
- (32) Goeman, J. J.; Van De Geer, S. A.; Van Houwelingen, H. C. *J. R. Stat. Soc. Series B Stat. Methodol.* **2006**, *68*, 477-493.

- (33) Burdinski Jr, T. K. , *Multiple Linear Regression Viewpoints*, **2000**, 2, 15-28
- (34) Healy, M. *Appl. Stat.* **1968**, 157-161.
- (35) Kessner, D.; Chambers, M.; Burke, R.; Agus, D.; Mallick, P. *Bioinformatics* **2008**, 24, 2534-2536.
- (36) Race, A. M.; Styles, I. B.; Bunch, J. J. *Proteomics* **2012**, 75, 5111-5112.
- (37) Mahalanobis, P. C. *Proc. Natl. Inst. Sci. (Calcutta)* **1936**, 2, 49-55.
- (38) Kendall, M. G. *A Course in the Geometry of n Dimensions*; Courier Corporation, 2004.
- (39) Lein, E. S.; Hawrylycz, M. J.; Ao, N.; Ayres, M.; Bensinger, A.; Bernard, A.; Boe, A. F.; Boguski, M. S.; Brockway, K. S.; Byrnes, E. J. *Nature* **2007**, 445, 168-176.
- (40) Ipsen, A. *Anal. Chem.* **2015**, 87, 1726-1734.
- (41) Du, P.; Stolovitzky, G.; Horvatovich, P.; Bischoff, R.; Lim, J.; Suits, F. *Bioinformatics* **2008**, 24, 1070-1077.
- (42) Palmer, A. D. *Information processing for mass spectrometry imaging*. Ph.D Thesis, University of Birmingham **2014**.
- (43) Xu, J. J. *Retrospective Theses and Dissertations* **1996**, 3120.

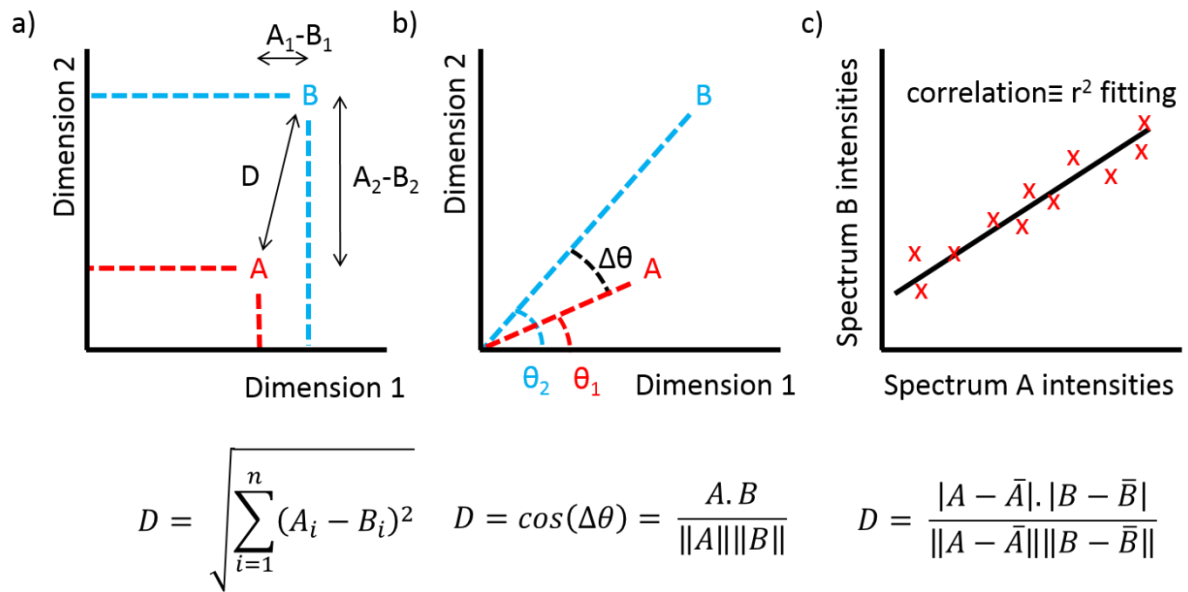


Figure 1. Visual representations of three of the distance metrics, a) Euclidean distance, b) cosine distance, and c) correlation.

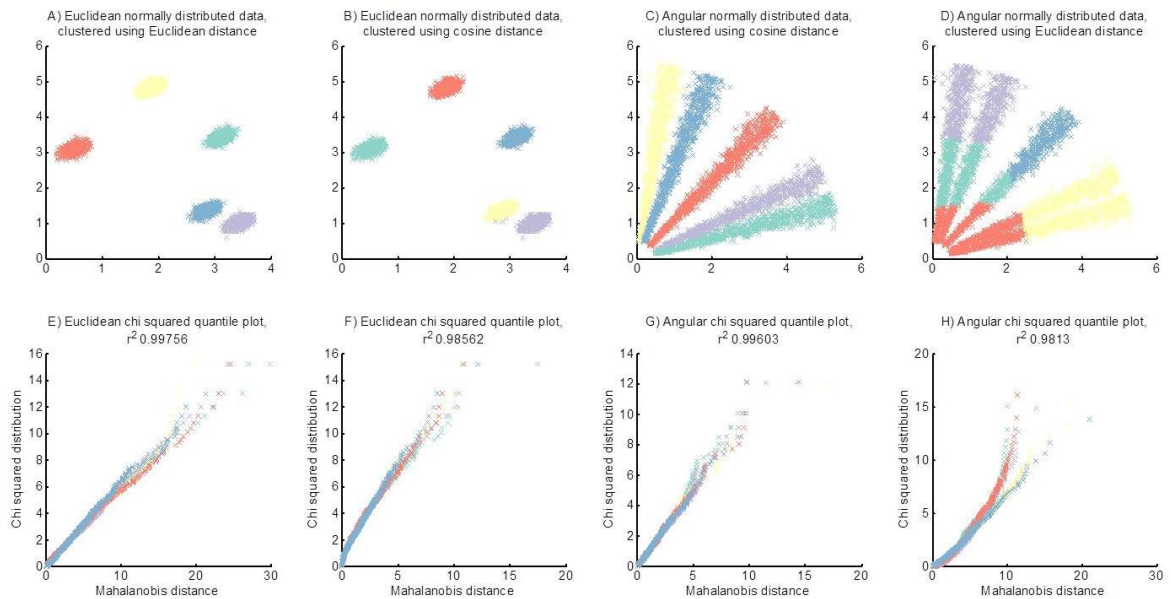


Figure 2. Simulated data of five clusters with normally distributed data (A, B), and angular normally distributed data (C, D), clustered using the Euclidean distance (A,D) and the cosine distance (B, C), with corresponding Chi squared quantile-normality plots below (E-H).

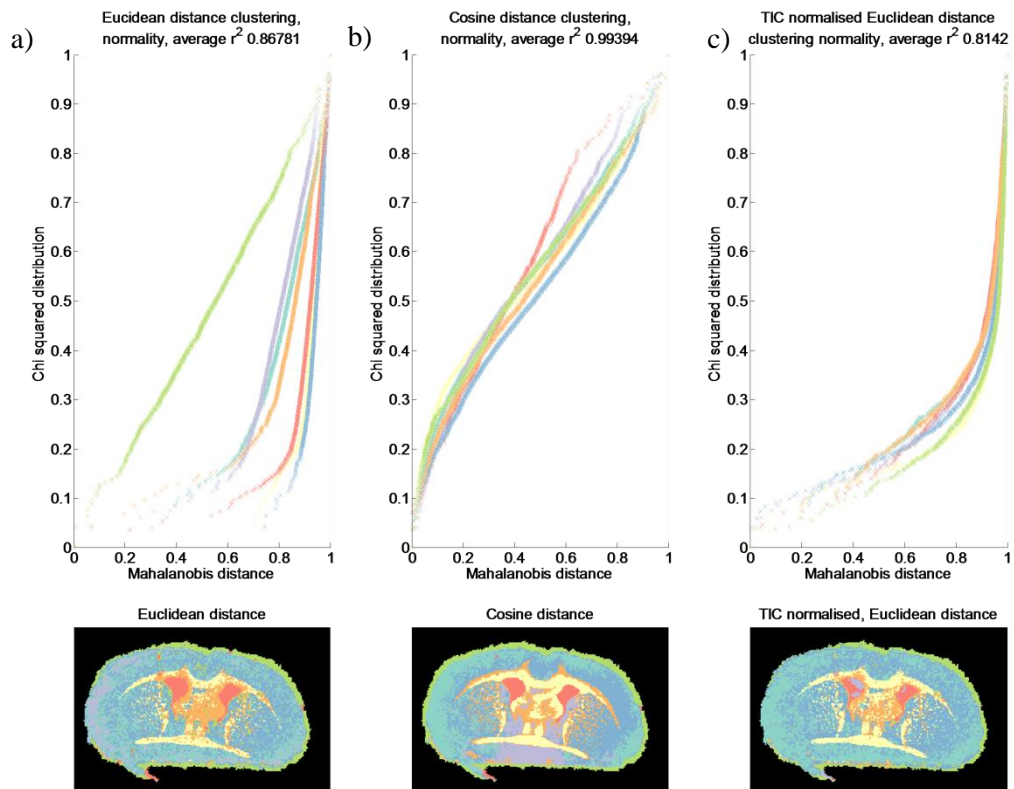


Figure 3. Quantile-Quantile plot in a) Euclidean space b)) angular space , and c) TIC normalised Euclidean space for the data within each of the 7 clusters of the coronal rat brain image segmented using a) Euclidean distance, b) cosine distance and c) Euclidean distance with TIC normalisation.

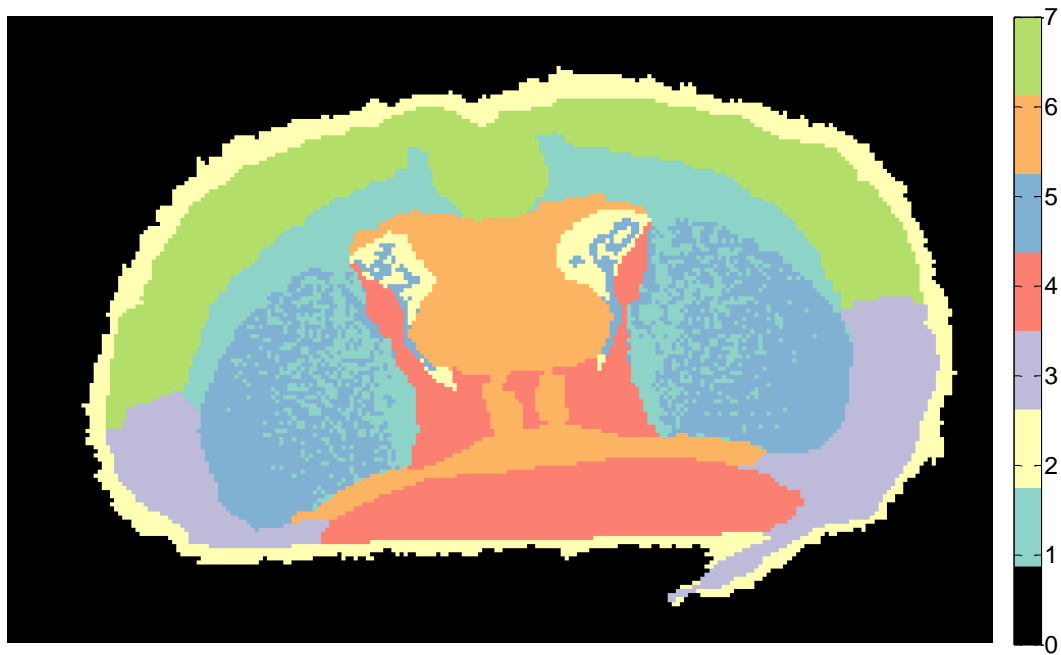


Figure 4. Seven anatomical regions used as reference data for statistical modelling, segmented based on a combination of comparison of a high resolution optical image (Figure S13) with selected ion images and PCA scores images along with comparison to the Allen brain atlas.

ID	Region	Polar Normality (r^2)	Euclidean normality (r^2)
1	Corpus callosum	0.983	0.822
2	Outer boundary	0.983	0.904
3	Olfactory areas	0.993	0.889
4	Brain stem	0.988	0.854
5	Caudoputamen	0.994	0.948
6	Lateral septal complex	0.984	0.916
7	Isocortex	0.990	0.680

Table 1. Results of r^2 fitting of the chi-squared quantile-quantile plots on the seven anatomical regions in Figure 4 showing that the data is closer to normal in polar space than in Euclidean.

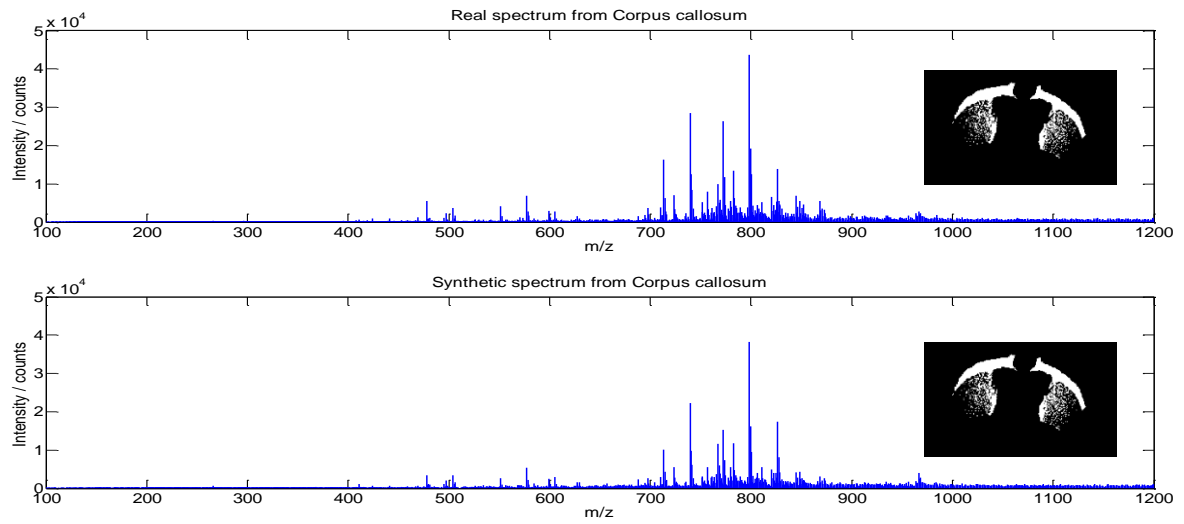


Figure 5. Top, real spectrum from the corpus callosum region of the reference data; bottom, synthetic data sampled from the multivariate normal distribution of the corpus callosum region, a high degree of spectral similarity is observed between the spectra

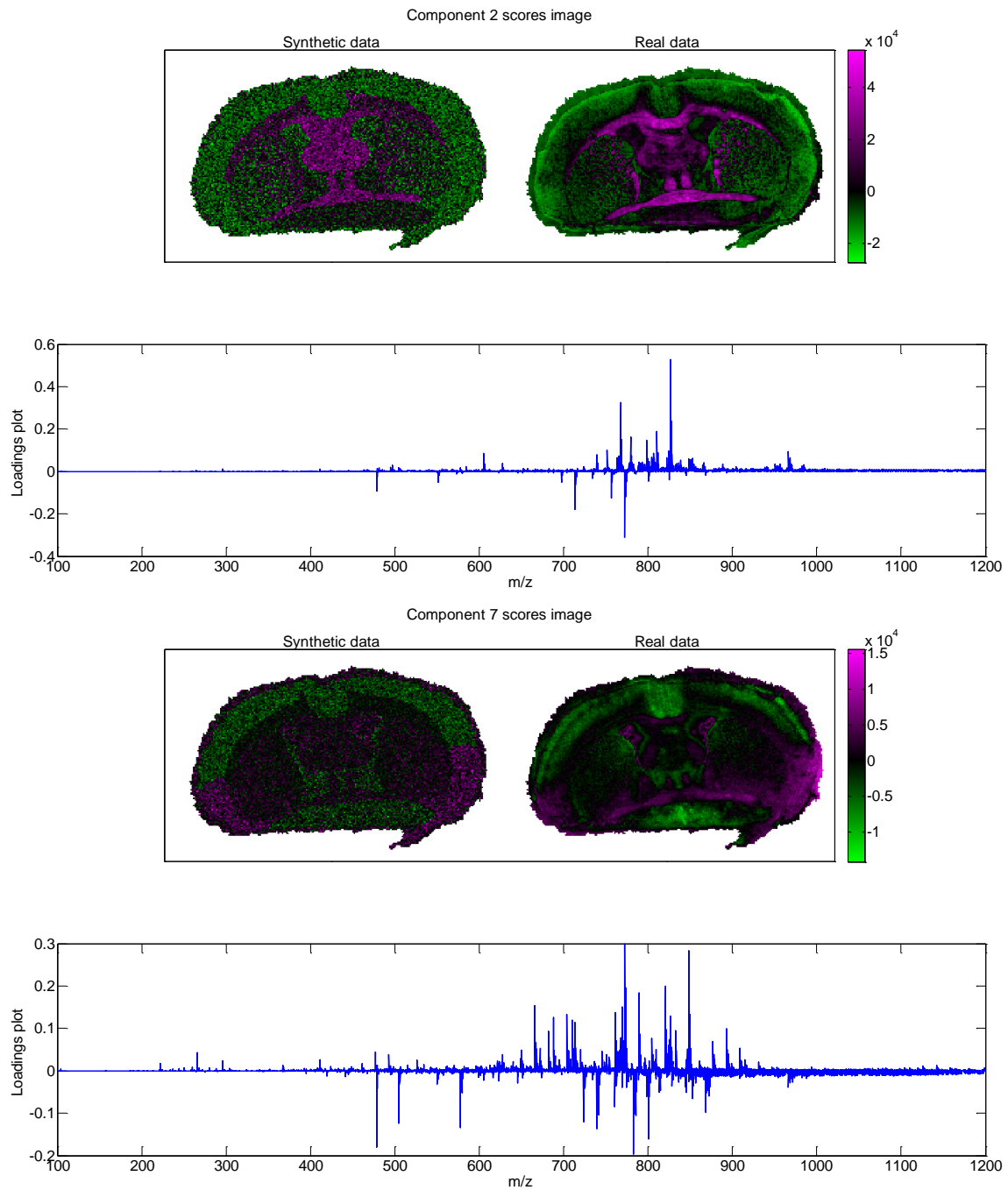


Figure 6. Results of PCA on the combined dataset containing both real and synthetic data, with scores images on top and projection loadings plots on the bottom. No principal component was found to separate the real from the synthetic data, indicating that any variance in the data is from the inherent biological and experimental variance in the reference data rather than introduced by the statistical modelling.

