

Using Shape Metrics to Describe 2D Data Points

William Franz Lamberti
Center for Public Health Genomics
Department of Biomedical Engineering
School of Data Science
University of Virginia
Charlottesville, VA, United States
william.f.lamberti@virginia.edu

Abstract—Traditional machine learning (ML) algorithms, such as multiple regression, require human analysts to make decisions on how to treat the data. These decisions can make the model building process subjective and difficult to replicate for those who did not build the model. Deep learning approaches benefit by allowing the model to learn what features are important once the human analyst builds the architecture. Thus, a method for automating certain human decisions for traditional ML modeling would help to improve the reproducibility and remove subjective aspects of the model building process. To that end, we propose to use shape metrics to describe 2D data to help make analyses more explainable and interpretable. The proposed approach provides a foundation to help automate various aspects of model building in an interpretable and explainable fashion. This is particularly important in applications in the medical community where the ‘right to explainability’ is crucial. We provide various simulated data sets ranging from probability distributions, functions, and model quality control checks (such as QQ-Plots and residual analyses from ordinary least squares) to showcase the breadth of this approach.

Index Terms—Shape Analysis, Image Measures, Explainability, Interpretability, Emerging Applications and Systems

I. INTRODUCTION

In the age of big data, 2D representations of multivariate data are still commonplace and a requirement for many experiments [1]–[5]. For example, analysts continue to simplify results to 2D representations by analyzing various underlying properties and implicit assumptions. However, many of these analyses are subjective and not numeric. For instance, analyzing the residuals of multiple regression is vital to ensure that the assumptions of the model are met [6], [7]. However, these checks are usually checked by human analysts and are not verified using numeric measures. There are many other aspects of multiple regression that require human decisions to be made such as the inclusion of interaction terms. These human decisions can be a downside when compared to deep learning approaches since deep learning models are able to learn what features are important for a given analysis without the need for a human analyst [8]. Having these human elements makes experiments using these approaches less explainable and interpretable for those not involved in the experiment.

Ensuring that a computational analysis is as explainable and interpretable as possible is key for explainable artificial intelligence (XAI) applications. XAI is used in various fields such as medicine and national security [9]. For example,

the European Union has passed laws ensuring a patient’s “right to explainability” [10]. In short, this law states that if a computational model is used to help make a diagnosis, the aspects of the computational model must be able to be described in layman’s terms. However, XAI is not limited to the intersection of science and policy. XAI is key for making scientific inferences since scientists need to understand how AI models use features to better understand the scientific phenomena. Thus, there is a fundamental need to make computational analyses as intuitive and clear as possible. A shape metric like area is a prime example of a metric which has a clear definition and is a concept that is understandable to the average person [9]. Thus, we desire to use explainable and interpretable shape metrics to describe data to help quantify the shapes of 2D data. This paper posits that all data that resides in 2D feature spaces can be represented as images. These images can be analyzed by extracting various useful shape metrics.

The work that provides the foundation for the idea of analyzing 2D data using shape metrics is eigenvalue decomposition. Eigenvalues correspond to the relative length of the axes of the data [9]. For example, if we observe a 2D scatterplot, the major and minor axis lengths will be captured by the first and second eigenvalues, respectively [9]. Thus, eigenvalue decomposition is the first idea that measures the shape of data. However, this idea has not been significantly expanded. To that end, we are providing foundational evidence and an approach to use a variety of shape metrics to describe data in 2D feature spaces.

We provide a new manner to analyze 2D data as images. We first convert the data in 2D space to 2D images of the shapes. We then collect shape metrics from the images. We lastly analyze the images for a given analysis. By quantifying 2D data using tangible shape metrics, we make analyses of 2D data more explainable and interpretable. An overview of our contribution is provided in Figure 1. The code for our experiments are provided at our GitHub link: https://github.com/billy1320/2d_shape_points.

II. METHODS AND MATERIALS

There are various simulated scenerios we will analyze: the discrimination between different Normal distributions, the detection of outliers in QQ-Plots, the discrimination of different 2D functions, and the analysis of multiple regression or ordinary least squares (OLS) residual analysis of variance.

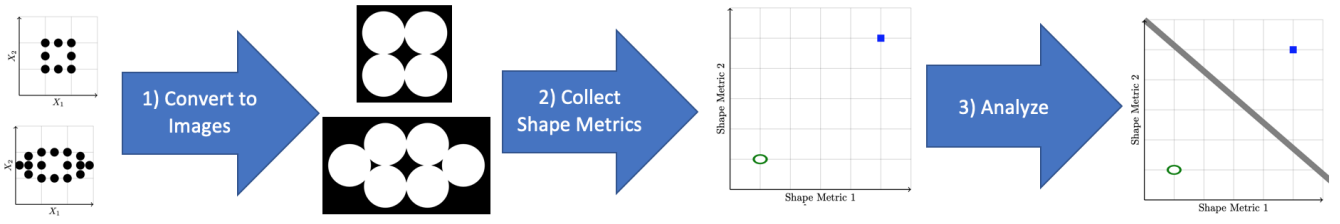


Fig. 1: Overview of analyzing 2D data as images using shape metrics. Once the 2D data is collected, step 1) converts the data to binary images. Step 2) collects various shape metrics of interest to describe the images. Step 3) analyzes the shape metrics. In this example, the differently shaped data points are classified. However, other analyses could be performed.

Normal Distributions: The discrimination of different Normal distributions will help to provide a foundation by which we are able to understand how shape metrics are useful on conceptually concrete examples. We will denote Normal or Gaussian distributions with a mean μ and variance σ^2 as $N(\mu, \sigma^2)$. The first experiment will have two simulated Normal distributions with a common variance, but different means. This experiment will show that the means of the distribution are not important for discriminating the distributions. The following two experiments will have pairs of Normals with the same mean but differing variances. These two experiments provide evidence that the shape metrics are useful for discriminating Normals with differing variances.

The followup to these experiments uses various mixtures of Gaussian distributions. This will show that different mixtures of Normals can be discriminated using shape metrics. Examples are provided in Figure 2. Thus, this set of four experiments using a variety of different Normal distributions shows that the variance and the number of Gaussians mixtures present are key for providing differently shaped data.

QQ-Plots and Outlier Detection: The second experiment aims to show that shape metrics can be used to classify QQ-Plots that have outliers. QQ-Plots are a visualization technique used to ensure that data follows a particular distribution [11]. For OLS, we want the residuals to follow a Normal distribution [6], [7], [12]. We simulated various QQ-Plots with no outliers (1000 random $N(0, 1)$), minor outliers (990 random $N(0, 1)$ and 10 random $N(3, 1)$), medium outliers (990 random $N(0, 1)$ and 10 random $N(5, 1)$), and major outliers (990 random $N(0, 1)$ and 10 random $N(10, 1)$) to showcase the ability of our approach to identify outliers. Examples of these QQ-Plots as images are provided in Figures 3a - 3d, respectively.

Functions: The third experiment aims to show that different shaped 2D functions can be discriminated using shape metrics. Each function had 1000 random simulated observations per resulting image. The first function was

$$Y = 3X + \epsilon \quad (1)$$

where $X \sim N(0, 100)$ and $\epsilon \sim N(0, 1)$. This is referred to as the “Linear” function. Note that \sim represents “distributed as”. The second function was

$$Y = 4 \sin X + \epsilon \quad (2)$$

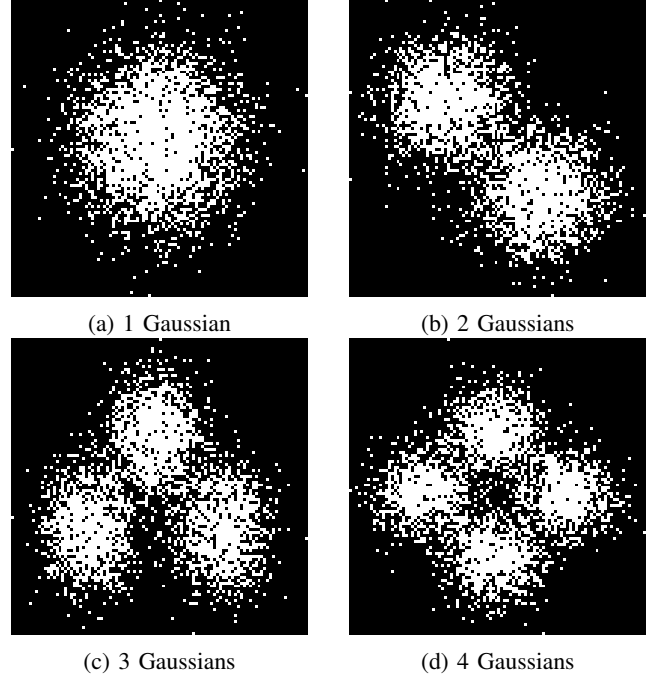


Fig. 2: Examples of mixtures of Gaussians.

where $X \sim N(0, 100)$ and $\epsilon \sim N(0, 0.25)$. This is referred to as the “Sine” model. The third function was

$$Y = X^2 + \epsilon \quad (3)$$

where $X \sim N(0, 100)$ and $\epsilon \sim N(0, 1)$. This is referred to as the “Parabola” function. The fourth function was

$$Y = X^4 + 10X^3 - 7X^2 + \epsilon \quad (4)$$

where $X \sim N(0, 100)$ and $\epsilon \sim N(0, 1)$. This is referred to as the “polynomial” or “Poly.” model. Examples of these functions are provided in Figures 4a - 4d, respectively.

OLS Residual Analysis: The fourth and last experiment corresponds to evaluating different types of variance plots for multiple regression. The plots should have random scatter and not display any obvious patterns. An example of this is provided in Figure 5a. Cone-like shape could indicate an increase in variance as the response increases. This is common for Poisson phenomena. Another concerning pattern

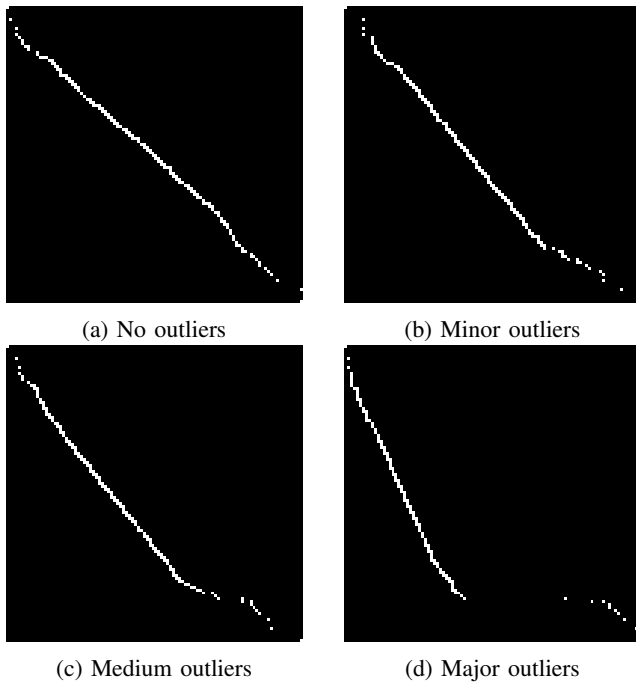


Fig. 3: QQ-Plot examples with different levels of outliers.

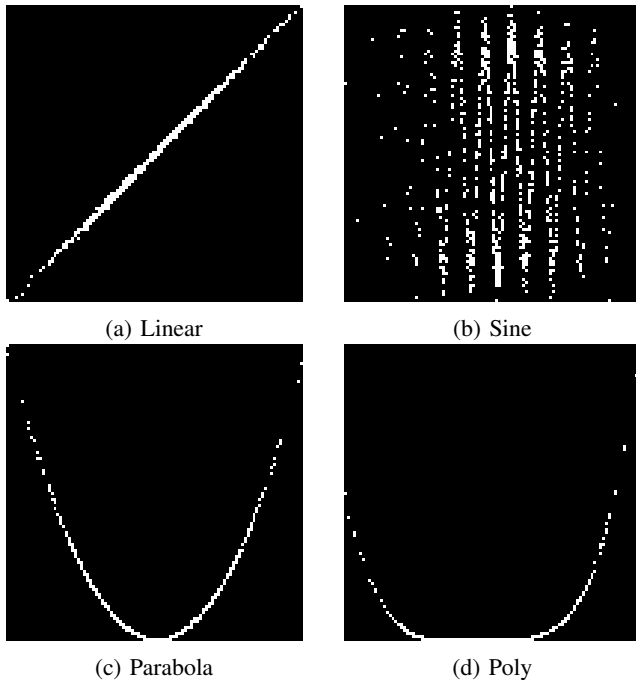


Fig. 4: Examples of the various simulated functional models resulting images.

looks like an almond or an eye. This occurs when the response is a proportion from a Binomial random variable (shorted to “Binom”). The last error discussed here is multiplicative errors (shorted to “Multi.”), where the plot will look like a bowtie. An example of multiplicative errors is provided in Figure 5b. Once these patterns are identified, transformations can be applied to

the response to correct for these errors [6]. Automating these transformations would help to make multiple regression more standardized and reproducible.

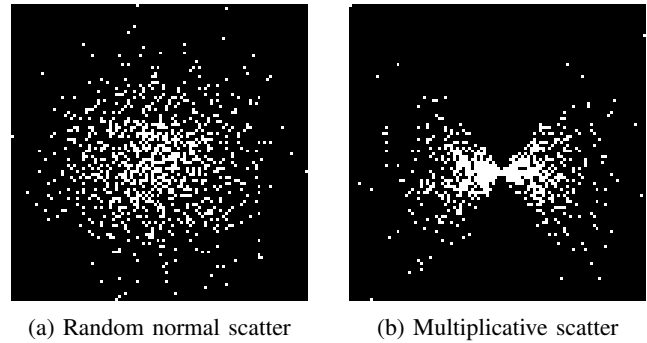


Fig. 5: Examples of a subset of the OLS residual simulated resulting images.

Converting Data to Images: The crucial step is converting the raw data into images. This is done by converting the images using 2D histograms and then accepting all positive signals. We will describe this more precisely using image operator notation [13]. Thus, 2D raw data are converted into 2D images using

$$\mathbf{b}[\vec{x}] = \Gamma_{>0} \mathbb{H}_{2,2} X, \quad (5)$$

where X is the input data, $\mathbb{H}_{2,2}$ converts the data into a 2D histogram [14], and $\Gamma_{>0}$ is the threshold image operator. From here, we would collect various shape metrics of interest.

Shape Metrics: Due to space limitations, we do not provide the exact manner in which these shape metrics are collected. However, we will provide a brief description of the shape metrics used in this analysis. Much of the description of these metrics is borrowed from Lamberti [15]. Lamberti’s extended descriptions on the shape metrics used in this analysis can be found at various sources [9], [15], [16]. The first metrics were the shape proportions (SP) and encircled image-histograms (EI), which were collected from the shape proportion and encircled image-histogram (SPEI) algorithm [16]. The EI is the black and white pixel counts of the shape after the shape is placed in the minimum encompassing circle and then the minimum encompassing square [16]. In other words, this is the area and the surrounding area of the shape. The SP value is the proportion of the area of the shape relative to the sum of the EI. The other shape metrics collected that were used in the model were the eigenvalues of the shapes, eccentricity [13], and circularity [13], [17]. The eigenvalues measure the major and minor axes of the shape. Eccentricity is the ratio of the major axis over the minor one. These are calculated using the 1st and 2nd eigenvalues of the shapes, respectively [13]. Circularity measures how circular a given shape is. This results in a total of 7 total metrics used during our analyses. The used metrics are summarized in Table I.

Modeling: Once we collected the 7 metrics, collected 100 simulated cases (which results in 100 images) for each class for a given experiment. We used a classification tree to

TABLE I: Table provides the metrics used in this analysis on a given image, i . The first column is the q^{th} metric, where $q \in \{1, 2, \dots, 7\}$. The last column provides the number of times the classification trees used each metric. These variables make our models interpretable and explainable [9].

$\bar{m}_{q,i}^*$	Metric	Counts
1	White EI	6
2	Black EI	0
3	SP value	3
4	Eccentricity	4
5	1 st Eigenvalue	2
6	2 nd Eigenvalue	0
7	Circularity	1

discriminate between the different classes [18]. We used 80% of the data as training and 20% as validation. We used stratified random sampling to ensure that the proportions were evenly split between the different classes. On the training data, we used 5-fold cross-validation to select the complexity parameter [3], [19].

III. RESULTS

A summary of the results are provided in Table II. Each set of experiments is analyzed in more detail in the following sections. As a whole, these experiments show that raw 2D data can be analyzed as images using their shape metrics.

TABLE II: Validation accuracy measures and confidence intervals (CIs) of classification trees for different experiments. The first 4 rows correspond to the comparison of multivariate Normal distributions under different conditions. The fifth row corresponds to the QQ-Plot experiment. The sixth row corresponds to the Function experiment. The seventh row corresponds to the OLS residual analysis experiment. Note that $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ and $\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$.

Experiment	Accuracy	Accuracy 95% CI
$N_1(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma), N_2(\begin{bmatrix} 10 \\ 10 \end{bmatrix}, \Sigma)$	0.525	(0.3613, 0.6849)
$N_1(\mu, \Sigma), N_2(\mu, \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix})$	1.00	(0.9119, 1.00)
$N_1(\mu, \Sigma), N_2(\mu, \begin{bmatrix} 0.001 & 0 \\ 0 & 0.001 \end{bmatrix})$	1.00	(0.9119, 1.00)
4 Gaussian Mixtures	0.975	(0.9126, 0.997)
QQ-Plots and Outlier Detection	0.95	(0.8769, 0.9862)
Functions	0.9375	(0.8601, 0.9794)
OLS Residual Analysis	0.9375	(0.8601, 0.9794)

Normal Distributions: This analysis shows that shape metrics are useful for only discriminating the shape of data, not the location of the data. The first experiment shows where

shape metrics will be unhelpful as only the means differ between the distributions. The shape of the distribution is primarily described by the variance of the distribution.

QQ-Plots and Outlier Detection: Despite having varying levels of outliers, we were able to discriminate between different levels of outliers on our QQ-Plots. This shows that the shape metrics are useful in quantifying different levels of outliers present in 2D data. This would help to quantify and detect outliers automatically. Thus, we have evidence that shape metrics can be used to detect the presence of outliers automatically.

Functions: Shape metrics are useful for classifying different kinds of 2D functions from one another. This is useful for evaluating a variety of different functional shapes. This can be used to help guide analysts on the kind of model to utilize for a given analysis. This experiment provides evidence for automatic function determination.

OLS Residual Analysis: Shape metrics are useful for determining the kind of variance pattern observed in multiple regression residual plots. This helps to remove the more subjective aspects of model evaluation and help to provide a standard and reproducible process for modeling choices. This provides evidence for automated OLS residual analysis using shape metrics.

Useful Features: Table I shows the number of times each of the metrics collected were used across all of the experiments in the classification trees. The White EI (or area) was used the most out of all of the chosen metrics. Eccentricity was the second most used metric. Thus, White EI and Eccentricity should be used in future analyses and applications at a minimum.

IV. CONCLUSIONS

Shape metrics are useful for describing a variety of different 2D data scenarios. This provides strong evidence that all 2D raw data are images and can be analyzed as such. Thus, we can standardize traditional analyses like multiple regression to help automate modeling decisions and make them more explainable and interpretable to those not involved in the analysis. This paper provides a strong foundation for analysts to analyze 2D data as images using shape metrics.

V. ACKNOWLEDGEMENTS

UVA Engineering Graduate Writing Lab Peer Review Group provided valuable feedback during initial drafts of this paper. We would also like to thank the Zang Lab for Computational Biology at the University of Virginia for their support.

REFERENCES

- [1] L. McInnes, J. Healy, and J. Melville, "UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction," *arXiv:1802.03426 [cs, stat]*, Sep. 2020, arXiv: 1802.03426. [Online]. Available: <http://arxiv.org/abs/1802.03426>
- [2] L. J. Mateo, S. E. Murphy, A. Hafner, I. S. Cinquini, C. A. Walker, and A. N. Boettiger, "Visualizing DNA folding and RNA in embryos at single-cell resolution," *Nature*, vol. 568, no. 7750, pp. 49–54, Apr. 2019. [Online]. Available: <https://www.nature.com/articles/s41586-019-1035-4>
- [3] Y. Takei, J. Yun, S. Zheng, N. Ollikainen, N. Pierson, J. White, S. Shah, J. Thomassie, S. Suo, C.-H. L. Eng, M. Guttman, G.-C. Yuan, and L. Cai, "Integrated spatial genomics reveals global architecture of single nuclei," *Nature*, vol. 590, no. 7845, pp. 344–350, Feb. 2021, number: 7845 Publisher: Nature Publishing Group. [Online]. Available: <https://www.nature.com/articles/s41586-020-03126-2>
- [4] T. Alexandrov, "Spatial Metabolomics and Imaging Mass Spectrometry in the Age of Artificial Intelligence," *Annual Review of Biomedical Data Science*, vol. 3, Jul. 2020. [Online]. Available: <https://papers.ssrn.com/abstract=3658948>
- [5] H. Hotelling, "Analysis of a complex of statistical variables into principal components," *Journal of Educational Psychology*, vol. 24, no. 7, pp. 498–520, Oct. 1933, publisher: Warwick & York. [Online]. Available: <http://mutex.gmu.edu/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=pdh&AN=1934-01354-001&site=ehost-live>
- [6] W. Mendenhall and T. T. Sincich, *A Second Course in Statistics: Regression Analysis*, 7th ed. Boston, MA: Pearson, Jan. 2011.
- [7] N. R. Draper and H. Smith, *Applied Regression Analysis*, third edition ed. New York: Wiley-Interscience, Apr. 1998.
- [8] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, G. Wang, J. Cai, and T. Chen, "Recent advances in convolutional neural networks," *Pattern Recognition*, vol. 77, pp. 354–377, May 2018. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0031320317304120>
- [9] W. F. Lamberti, "An Overview of Explainable and Interpretable Artificial Intelligence," in *AI Assurance: Towards Valid, Explainable, Fair, and Ethical AI*. Elsevier, 2022.
- [10] E. Union, "Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation) (Text with EEA relevance)," May 2016, legislative Body: EP, CONSIL. [Online]. Available: <http://data.europa.eu/eli/reg/2016/679/oj/eng>
- [11] M. B. Wilk and R. Gnanadesikan, "Probability Plotting Methods for the Analysis of Data," *Biometrika*, vol. 55, no. 1, pp. 1–17, 1968. [Online]. Available: <https://www.jstor.org/stable/2334448>
- [12] G. K. Bhattacharyya and R. A. Johnson, *Statistical Concepts and Methods*, 1st ed. Wiley, 1977. [Online]. Available: <https://www.wiley.com/en-us/Statistical+Concepts+and+Methods-p-9780471072041>
- [13] J. M. Kinser, *Image Operators: Image Processing in Python*, 1st ed. Boca Raton, FL: CRC Press, Oct. 2018.
- [14] "numpy.histogram2d — NumPy v1.21 Manual." [Online]. Available: <https://numpy.org/doc/stable/reference/generated/numpy.histogram2d.html>
- [15] W. F. Lamberti, "Classification of Synthetic Aperture Radar Images of Icebergs and Ships Using Random Forests Outperforms Convolutional Neural Networks," in *2020 IEEE Radar Conference (RadarConf20)*, Sep. 2020, pp. 1–6, iSSN: 2375-5318.
- [16] —, "Algorithms to Improve Analysis and Classification for Small Data," Ph.D., George Mason University, United States – Virginia, 2020, iSBN: 9798557033350. [Online]. Available: <http://search.proquest.com/docview/2476825035/abstract/90CC4207B46B4068PQ/1>
- [17] A. Rosenfeld, "Compact Figures in Digital Pictures," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. SMC-4, no. 2, pp. 221–223, Mar. 1974.
- [18] T. Therneau, B. Atkinson, B. R. p. o. t. i. R. port, and maintainer 1999-2017), "rpart: Recursive Partitioning and Regression Trees," Apr. 2019. [Online]. Available: <https://CRAN.R-project.org/package=rpart>
- [19] G. James, D. Witten, T. Hastie, and R. Tibshirani, Eds., *An introduction to statistical learning: with applications in R*, ser. Springer texts in statistics. New York: Springer, 2013, no. 103, oCLC: ocn828488009.