1 **TITLE**

2 Multidimensional scaling informed by *F*-ratio for hypothesis testing in microbial community analysis

3 **AUTHORS AND AFFILIATIONS**

4 Hyungseok Kim[*,1,2], Soobin Kim[*,3], Megan M. Morris[4], Jeff A. Kimbrel[4], Xavier Mayali[4], and Cullen R.

5 Buie[1]

6 [1] Department of Mechanical Engineering, Massachusetts Institute of Technology, Cambridge MA USA

7 [2] Institute for Data, Systems, and Society, Massachusetts Institute of Technology, Cambridge MA USA

8 [3] Department of Statistics, University of California, Davis, Davis CA USA

9 [4] Physical and Life Sciences Directorate, Lawrence Livermore National Laboratory, Livermore CA USA

10 [*] Equal contribution

11 **SUMMARY**

12 Multidimensional scaling (MDS) is a common dimension reduction method to capture a global pattern

13 and structure in microbial community. While a confirmatory analysis with the classical MDS is a

14 promising approach in broadening an interpretation of the multivariate data (e.g., building a classifier,

15 hypothesis testing), little attention is gained because the configuration can be biased towards the external

16 information. Here we propose a confirmatory MDS informed by an *F*-ratio addressing a community-

17 structural hypothesis testing result under a binary setting. Using a simulated or a 16S rRNA microbial

18 community example, we present how the proposed configuration incorporates the testing result. We

19 evaluate a performance by comparing to the previous confirmatory MDS, demonstrating that our method

20 is less dependent on a selection of a model hyperparameter and minimally alters the classical MDS

21 configuration. Our method proposes a broader applicability of MDS in modern biological and

22 multivariate data analysis.

23 **INTRODUCTION**

24    Biotechnological advances in the past several decades have expanded a size and features of the

25    multivariate data, necessitating a dimensionality reduction as a tool for the interpretation. By extracting an

26    essential information from the biological data, the dimensionality reduction seeks a visual representation

27    of the multivariate in a lower dimensional space. To retain a consistent data structure while performing

28    the dimension reduction, a configuration is sought in a way that preserves dispersion or dissimilarity

29    between samples, a process called as the multidimensional scaling (MDS). Compared to other nonlinear

30    methods in dimensionality reduction [1-3], MDS is known to retain a global structure and represent a

31    long-range interaction between samples, allowing its popularity for a long time since its inception.

32    In a classical MDS, the configuration is determined in a way that minimizes the difference between the

33    dissimilarity in the original and the low dimensional space, a measure termed the stress. In microbial

34    ecology community where the input data are compositional (i.e., a sample is comprised of different

35    species expressing a level of abundance), the dissimilarity is measured by a difference in expression level

36    of species within samples. For example, when interpreting a microbial community dataset such as 16S

37    rRNA gene expression, a distance metric so-called the Unifrac [4] is used, and it allows to incorporate

38    phylogenic diversity as well as the compositional differences. The choice of an appropriate dissimilarity

39    metric remains an important criterion for processing the input data and represent them in the low

40    dimensional space.

41    In addition to visualizing the community structure with the dimensionality reduction, a quantitative

42    analysis is carried out using the statistical inference such as hypothesis testing. Because the compositional

43    abundances do not assume a probabilistic distribution *a priori* (highly skewed, zero-inflated [5]), a

44    nonparametric model is preferred over parametric models for constructing a statistic. When the response

45    is assumed to be independent between the samples, the statistic can be readily obtained by permuting the

46    data labels [6]. For testing a difference in the ecological groups based on the sample dispersion, a

47    (pseudo) *F* statistic has widely been used for each permutation [5, 7].

48 While hypothesis testing provides a quantitative perspective to understand the multivariate data, it should

49 be noted that the testing result does not account for the MDS configuration. This is because most MDS

50 (e.g., PCA, PCoA) do not consider data labels, whereas the hypothesis testing aims to infer whether

51 dependent variables (e.g., class, label) are influenced by distribution of each sample group. For example,

52 a configuration from the classical MDS is not able to explain a small but statistically meaningful

53 difference between groups of different treatment.

54 The insufficient explanation by the classical MDS encourages to revise the approach and to address a

55 structural hypothesis by including an external information, conferred by the responses or class [8] labels.

56 Because the stress function is non-convex and an optimization algorithm can produce several local

57 configurations, there is a notion allowing an altered configuration up to a point where its stress does not

58 deviate too much from the classical configuration [9]. Broadly termed as the confirmatory MDS, it

59 imposes an external constraint to the classical MDS carrying over an additional task to minimizing the

60 stress.

61 In the recent confirmatory MDS methods, an objective function is constructed by adding a confirmatory

62 term to the stress, former of which can be quantified by the labels then multiplied by a hyperparameter.

63 While these confirmatory MDS methods have successfully visualized the multivariate structure in a way

64 that differentiates each sample group with a discriminative purpose [10, 11], choosing a proper

65 hyperparameter remains as a bottleneck towards a broader application of these inventive methods. For

66 example, setting a high hyperparameter results in an undesirable stress and a misleading configuration

67 distorted from the original.

68 In this study we propose an alternative MDS informed by a hypothesis testing inferences under a binary

69 class setting. Our approach is motivated by a purpose to explain a statistical difference between groups, if

70 any, using the MDS combined with a confirmatory analysis. Because the method does not target to

71 directly discriminate between groups, the motivation distinguishes itself from previously proposed

72   confirmatory MDS. Furthermore, by characterizing the proposed framework we show that the revised

73   configuration is less dependent to the choice of the model hyperparameter, mediating the previous issue

74   with the distortion.

## METHODS

### Problem formulation

77   Consider a balanced design where the number of total observations is $N$, and each observation $x_i$ is $S$-

78   dimensional, pertaining to a set of labels $y_i \in \{0,1\}$ for every $i = 1 \dots N$. Based on the observations $(x_1, \dots$

79   $x_N)$, a distance matrix is obtained as $\mathbf{d} = [d_{ij}] \in \mathrm{R}^{N \times N}$. Now given the distance $\mathbf{d}$, we seek a two-

80   dimensional configuration $\mathbf{z} = (z_1, \dots z_N) \in \mathrm{R}^{N \times 2}$, that best represents the original dimension by the

81   following criteria.

### Classical MDS

83   In classical MDS, a configuration is realized by minimizing the following objective function.

$$O(\mathbf{z}) = \frac{1}{2} \sum_{i,j} (d_{ij} - \|\mathbf{z_i} - \mathbf{z_j}\|_2)^2 \tag{1}$$

84

85   In other words, the configuration $\mathbf{z}$ is obtained in a way that tries to preserve a distance between a pair of

86   observations $(x_i, x_j)$ for each $i, j \in N$. Note that Equation (1) does not contain any terms related to $y_i$,

87   meaning that classical MDS does not consider the class labels.

### Confirmatory MDS for F-informed hypothesis testing

89   *Hypothesis testing for non-parametric multivariate analysis of variance*

90   When testing a statistical difference between groups in multivariate analysis, the group variance is the

91   measure of interest. It is represented by the *F*-statistic, a ratio between two group variances, each

92    respectively derived from across- and within-group. While the conventional *F*-test requires an assumption

93    that each observation follows a normal distribution, the convention is generalized by introducing an

94    analogous statistic (pseudo *F*-ratio) which is combined with label permutation for quantifying a statistical

95    significance. In this non-parametric approach [7], the pseudo *F*-ratio is defined as

$$F = \frac{\sum_{i,j} d_{i,j}^2 - 2\sum_{i,j} \mathbb{1}\{y_i = y_j\} d_{ij}^2}{2\sum_{i,j} \mathbb{1}\{y_i = y_j\} d_{ij}^2} \cdot (N - 2),$$ (2)

96

97    where $\mathbb{1}\{\cdot\}$ denotes an indicator function. Since the pseudo *F*-ratio does not follow an *F*-distribution

98    under the relaxed model assumption, it is instead evaluated by an empirical distribution that is created by

99    'permuting' the labels. That is, in every permutation a new *F*-ratio, $F^\pi$, is obtained from the data structure,

100    and by repeating this, we have a *P*-value written as

$$P = \frac{\text{Number of case where } (F^\Pi \geq F)}{\text{Number of total repeat}}.$$ (3)

101

102    Known as the permutational multivariate analysis of variance (PERMANOVA), the hypothesis testing

103    method has a broad application in microbial community analysis.

104    *Proposed MDS*

105    Now we propose a new multidimensional scaling that incorporates a hypothesis testing result in the

106    multivariate setting. This is enabled by adding a confirmatory term to the classical MDS (Eq. (1)), giving

107    an objective function as

$$O(\mathbf{z}) = \underbrace{\frac{1}{2}\sum_{i,j}(d_{ij} - \|\mathbf{z}_i - \mathbf{z}_j\|_2)^2}_{\text{MDS term}} + \lambda \cdot \underbrace{\frac{1}{2}\left|\sum_{i,j}[1 - (f_{\mathbf{z}}(\Phi_o) + 1)\mathbb{1}\{y_i = y_j\}]\|\mathbf{z}_i - \mathbf{z}_j\|_2^2\right|}_{\text{Confirmatory term}}$$ (4)

108

109    where $\Phi_o$ is a ratio constant expressed in terms of the distance $\mathbf{d}$ and labels $y_i$ and is defined as

$$\Phi_o(\mathbf{d}, y) := \frac{\sum_{i,j} \mathbb{1}\{y_i \neq y_j\} d_{ij}^2}{\sum_{i,j} \mathbb{1}\{y_i = y_j\} d_{ij}^2},$$

110

111      and $f_{\mathbf{z}}(\Phi_o) : \mathrm{R} \to \mathrm{R}$ is a mapping function which is determined by the configuration $\mathbf{z}$. An exact derivation

112      of $f_{\mathbf{z}}$ and a detailed description on the confirmatory term in Equation (4) is described in Appendix B.

113      Given Equation (4), we want to find an optimal configuration $\mathbf{z}^*$ such that $\mathbf{z}^* = \mathrm{argmin}_{(z1,\cdots,zN)} O(\mathbf{z})$.

114      *Algorithm for the proposed MDS*

115      Because the confirmatory term in (4) is strictly convex in terms of $\mathbf{z}$, we are able to minimize $O(\mathbf{z})$ by

116      using the Majorize-Minimization (MM) algorithm, a typical approach in the MDS optimization task [12].

117      While the implementation of MM algorithm is described in detail in Appendix C, we provide its update

118      rule as below.

---

**Algorithm 1** MM algorithm for pseudo $F$-informed MDS

For epoch $t$ and every $i = 1, \cdots N$,

$$\mathbf{z}_i^{[t+1]} \leftarrow \frac{2}{2(N-1) + \lambda\delta(\mathbf{z}_i^{[t]})(N - (N-2)f_{\mathbf{z}}(\Phi))}$$

$$\times \left[ (1 + \lambda\delta(\mathbf{z}_i^{[t]})) \sum_{\substack{j=1 \\ \epsilon_{ij}=0}}^{N} \mathbf{z}_j^{[t]} + (1 - \lambda f_{\mathbf{z}}(\Phi)\delta(\mathbf{z}_i^{[t]})) \sum_{\substack{j=1 \\ \epsilon_{ij}=1}}^{N} \mathbf{z}_j^{[t]} + \sum_{j=1}^{N} d_{ij} \frac{\mathbf{z}_i^{[t]} - \mathbf{z}_j^{[t]}}{\|\mathbf{z}_i^{[t]} - \mathbf{z}_j^{[t]}\|_2} \right],$$

where $\epsilon_{ij} = \mathbb{1}\{y_i = y_j\}$, $\delta_i(\mathbf{z}) = \mathrm{sign} \sum_{j=1}^{N}[1 - (f_{\mathbf{z}}(\Phi)+1)\epsilon_{ij}]\|\mathbf{z}_i - \mathbf{z}_j\|_2^2$, with an initial value obtained from a classical MDS.

---

119

## RESULTS AND DISCUSSION

120

121      Using Algorithm 1, we sought to determine how well our proposed MDS approach can produce a two-

122      dimensional configuration in simulated and experimental data. We then assess its performance by

123      comparing ours to a recent method in confirmatory MDS [13].

124      **Simulated data**

125      We first provide a representing case where our method can be useful in visualizing multidimensional data.

126 To do this we consider a binary labeled dataset where each group originates in a different multivariate

127 Gaussian distribution. In a three-dimensional setting, for example, consider a balanced design where an

128 observation expressed as

$$
x_i \sim \begin{cases} \mathcal{N}\left([0,0,0]^\top, \begin{bmatrix} 3 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 1 \end{bmatrix}\right), & i = 1, 2, \cdots 50 \\ \mathcal{N}\left([0,0,1]^\top, \begin{bmatrix} 3 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 1 \end{bmatrix}\right), & i = 51, 52, \cdots 100. \end{cases} \tag{5}
$$

129

130 As expected, PERMANOVA testing result indicates there is a statistically significant difference between

131 the groups with pseudo $F = 5.402$ and $p = 0.005$. However, a classical MDS does not distinguish groups

132 in two-dimensional configuration (Figure 1a) with $p = 0.914$, because the difference is in the third

133 dimension with the lowest variance among the principal diagonals (Figure 6 in Appendix C).

134 On the other hand, our MDS configuration is able to display the difference between the groups, and it

135 becomes clearer when a hyperparameter $\lambda$ is large (Figure 1b,c). The distinctions are also verified by a

136 low p-value resulting from PERMANOVA test using the two-dimensional configurations.
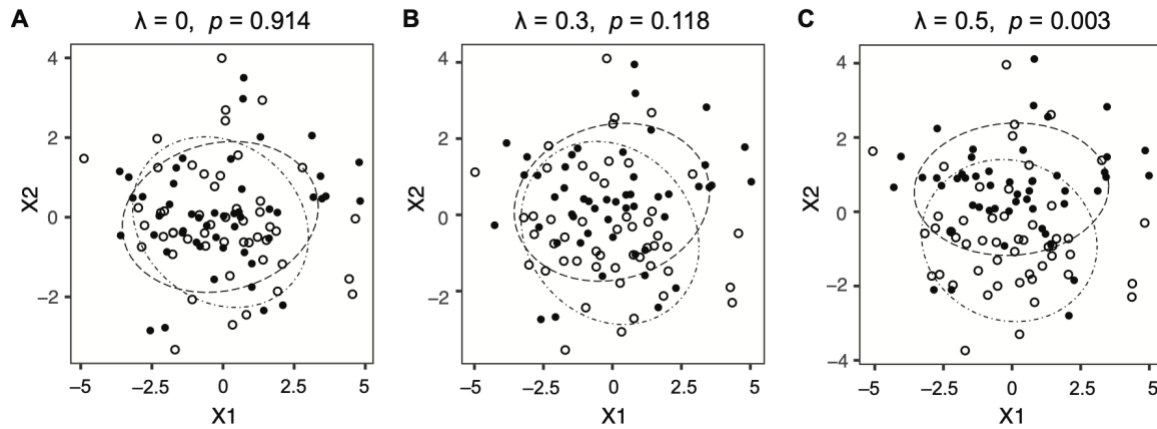


137

138    Figure 1: Two-dimensional visualization of proposed MDS with a hyperparameter (a) $\lambda = 0$ (classical

139    MDS), (b) $\lambda = 0.3$, and (c) $\lambda = 0.5$. For each configuration, a p-value is given based on PERMANOVA

140    test. An ellipse is drawn for each group with a confidence interval of 80%.

141    We next evaluate the performance of our method to existing confirmatory MDS by calculating a stress,

142    which is defined as

$$\text{Stress} = \frac{\sum_{i,j}(d_{ij} - \|\mathbf{z}_i - \mathbf{z}_j\|_2)^2}{\sum_{i,j} d_{ij}^2},$$

143

144    or by calculating a Spearman correlation from Shepard diagram, as measures of evaluating the

145    performance of MDS [14]. As a result, our proposed MDS visualization presents a Stress ~ 0.2 regardless

146    of a choice of a hyperparameter $\lambda$ (Figure 2a), suggesting either configuration can be used for visualizing

147    the simulated data [15]. This is in contrast to the previous method [13] where the stress increases as $\lambda$

148    becomes large, implying their approach distinguishes groups at the expense of the original distance

149    structure. Similarly, in Shepard plot with a choice of $\lambda$, our method presents a more consistent correlation

150    of the sample pair distance in between three- and two-dimension (Figure 2b,2c, and 7 in Appendix A).
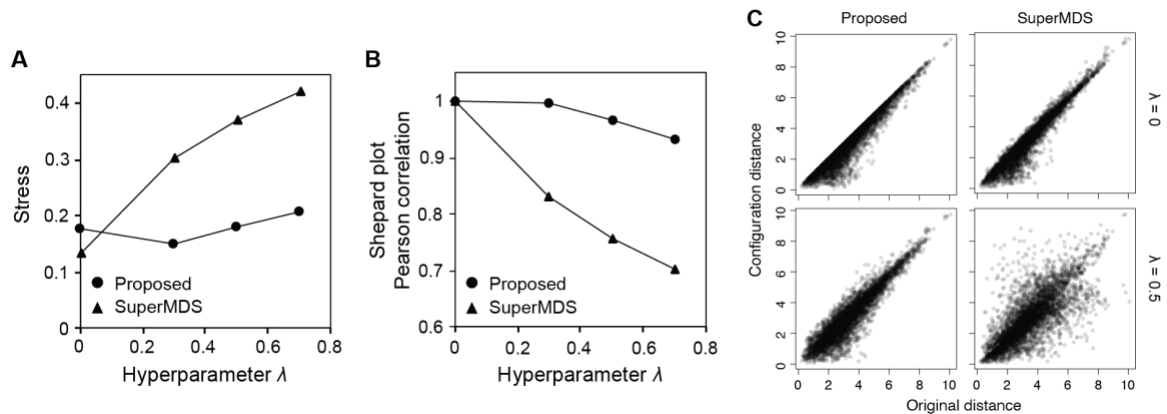


151

152    Figure 2: Performance of proposed MDS compared to an existing label-informed confirmatory MDS

153    (SuperMDS [13]) by using (a) Stress and (b) Spearman correlation from the simulated data. (c) Shepard

154  plot of the proposed MDS comparing to SuperMDS for a hyperparameter $\lambda = 0.5$. More results are

155  displayed in Figure 7.

156  It is worth noting that proposed MDS is invariant by the choice of the hyperparameter, which has not

157  been observed in the existing confirmatory MDS methods.

**Microbial community dataset**

159  Next, we provide an example where biological hypothesis testing result is conveyed to the MDS

160  configuration. We take microbial community dataset containing thirty-six, balanced samples of a binary

161  label (e.g., with or without a presence of microbial host) [16]. In detail, each data represents expression

162  levels 16S rRNA gene of 72 bacterial taxa, and the distance between samples is measured using the

163  weighted Unifrac [17]. Particular attention is made on these datasets, that the classical two-dimensional

164  MDS configuration does not explain PERMANOVA test results on a group differences. As shown in

165  Figure 3a, groups in site 1 are dispersed in a different location whereas site 2 groups are not, when

166  visualized using the classical MDS. In both sites, however, moderately small P-values are obtained ($< 0.1$,

167  Figure 3b), indicating the group difference in the community structure is, in fact, statistically significant.
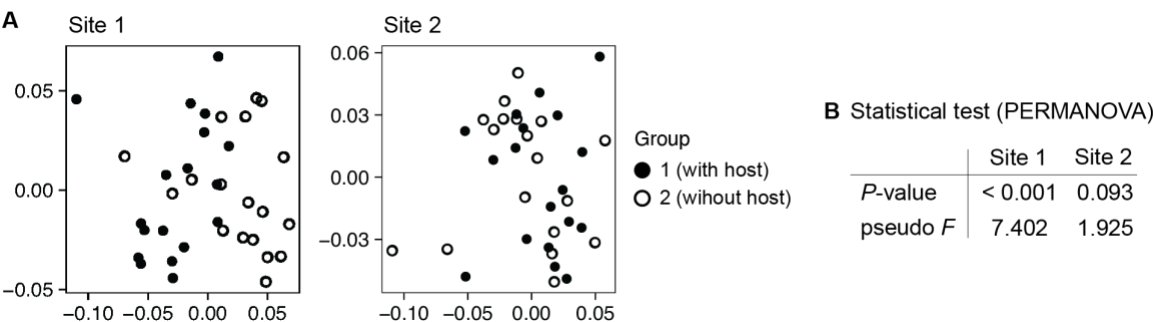


168

169  Figure 3: (a) Multidimensional scaling and (b) statistical test result on a group difference between two

170  sample groups for each site. The sample presents a microbial community mea- sured by 16S rRNA gene

171  expression.

172 Using the community dataset, we present a configuration with the proposed MDS visualization. As

173 expected, for site 1 community samples the configuration retains its distinction between the class labels

174 regardless of the choice of the hyperparameter $\lambda$ (Figure 4a-c). Moreover, for site 2 samples we observe a

175 higher distinction between the groups with increasing $\lambda$ (Figure 4d-f). The observation with the

176 visualizations is justified by a quantitative measure using P -value calculated on the 2D configurations

177 (Figure 4g).

178



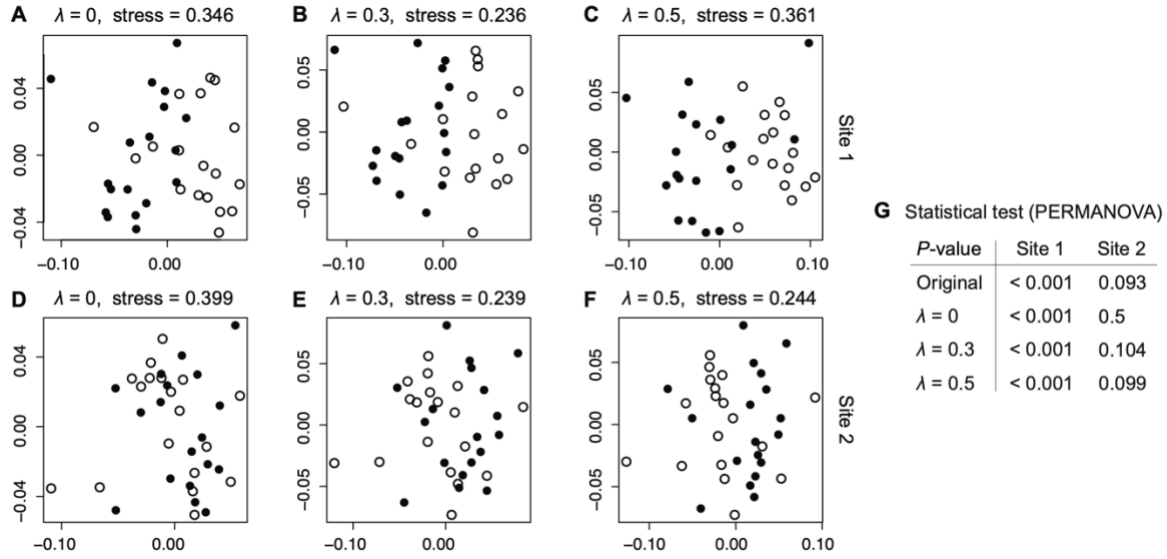| G Statistical test (PERMANOVA) | | |
| --- | --- | --- |
| P-value | Site 1 | Site 2 |
| Original | < 0.001 | 0.093 |
| $\lambda = 0$ | < 0.001 | 0.5 |
| $\lambda = 0.3$ | < 0.001 | 0.104 |
| $\lambda = 0.5$ | < 0.001 | 0.099 |

179 Figure 4: Two-dimensional configuration of microbial community samples using the proposed MDS

180 method, where samples are collected from (a-c) site 1 and (d-f) site 2. (g) Statistical significance on the

181 group difference between two treatments using PERMANOVA test.

182 We then evaluate the performance of the proposed MDS using stress measurement and Shepard plot.

183 Again, we observe stress does not strictly depend on the hyperparameter $\lambda$ or even show a decreased

184 value when $\lambda$ is nonzero (e.g., 0.1, 0.3) compared to the classical MDS ($\lambda = 0$, Figure 5a). Shepard plot

185 and Spearman correlation also show that the configurations nicely preserve the original distance in the

186 microbial community data, except for a case when the largest $\lambda$ is set to site 1 community.
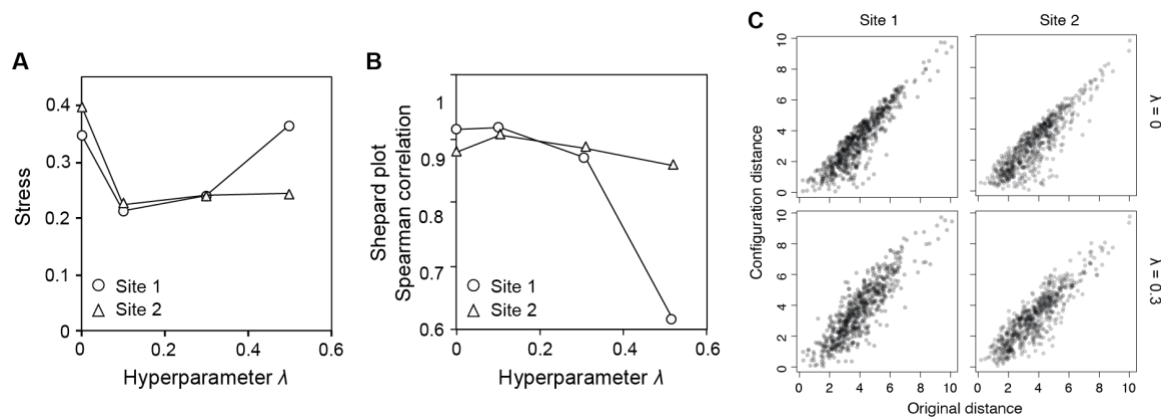
Figure 5: Evaluation of the proposed MDS using microbial community data, measured by (a) stress, (b) Spearman correlation, and (c) Shepard plot from each sample site. More results are displayed in Figure 8 (Appendix A).

**CONCLUSION**

- A new multidimensional scaling method which incorporates an F-statistic-informed hypothesis testing is proposed.

- We find that the performance of the proposed MDS excels existing MDS methods for addressing class labels, as evaluated by its stress and Shepard plot, validated using both simulated and real datasets.

- The proposed MDS can be useful when analyzing a biological dataset with F -informed hypothesis testing, providing informative and precise dimension reduction, especially for visualization.

- The method is less dependent on the choice of hyperparameter when producing the configuration. This lessens the risk of overreliance on the class labels in that the data are automatically grouped to a suitable degree. Also, users may avoid the hassle of hyperparameter selection using such as cross-validation.

**REFERENCES**

1.    Tenenbaum JB, Silva Vd, Langford JC: **A Global Geometric Framework for Nonlinear Dimensionality Reduction**. *Science* 2000, **290**(5500):2319-2323.

206    2.     McInnes L, Healy J, Melville J: **UMAP: Uniform Manifold Approximation and Projection for**
207            **Dimension Reduction**. In.; 2018: arXiv:1802.03426.

208    3.     van der Maaten L, Hinton G: **Visualizing Data using t-SNE**. *Journal of Machine Learning*
209            *Research* 2008, **9**(86):2579-2605.

210    4.     Lozupone C, Knight R: **UniFrac: a new phylogenetic method for comparing microbial**
211            **communities**. *Appl Environ Microbiol* 2005, **71**(12):8228-8235.

212    5.     Gijbels I, Omelka M: **Testing for Homogeneity of Multivariate Dispersions Using**
213            **Dissimilarity Measures**. *Biometrics* 2013, **69**(1):137-145.

214    6.     Holmes AP, Blair RC, Watson JD, Ford I: **Nonparametric analysis of statistic images from**
215            **functional mapping experiments**. *J Cereb Blood Flow Metab* 1996, **16**(1):7-22.

216    7.     Anderson MJ: **A new method for non-parametric multivariate analysis of variance**. *Austral*
217            *Ecology* 2001, **26**(1):32-46.

218    8.     Ding CS: **Testing Pattern Hypotheses with MDS**. In: *Fundamentals of Applied*
219            *Multidimensional Scaling for Educational and Psychological Research.* Edited by Ding CS.
220            Cham: Springer International Publishing; 2018: 165-173.

221    9.     Borg I, Groenen P: **Confirmatory MDS**. In: *Modern Multidimensional Scaling: Theory and*
222            *Applications.* Edited by Borg I, Groenen P. New York, NY: Springer New York; 1997: 181-197.

223    10.    Cox TF, Ferry G: **Discriminant analysis using non-metric multidimensional scaling**. *Pattern*
224            *Recognition* 1993, **26**(1):145-153.

225    11.    Yang F, Yang W, Gao R, Liao Q: **Discriminative Multidimensional Scaling for Low-**
226            **Resolution Face Recognition**. *IEEE Signal Processing Letters* 2018, **25**(3):388-392.

227    12.    Borg I, Groenen P: **A Majorization Algorithm for Solving MDS**. In: *Modern Multidimensional*
228            *Scaling: Theory and Applications.* Edited by Borg I, Groenen PJF. New York, NY: Springer New
229            York; 1997: 169-197.

230    13.    Witten DM, Tibshirani R: **Supervised multidimensional scaling for visualization,**
231            **classification, and bipartite ranking**. *Computational Statistics & Data Analysis* 2011,
232            **55**(1):789-801.

233    14.    Dexter E, Rollwagen-Bollens G, Bollens SM: **The trouble with stress: A flexible method for**
234            **the evaluation of nonmetric multidimensional scaling**. *Limnology and Oceanography:*
235            *Methods* 2018, **16**(7):434-443.

236    15.    Kruskal JB: **Multidimensional scaling by optimizing goodness of fit to a nonmetric**
237            **hypothesis**. *Psychometrika* 1964, **29**(1):1-27.

238    16.    Kim H, Kimbrel JA, Vaiana CA, Wollard JR, Mayali X, Buie CR: **Bacterial response to spatial**

239        **gradients of algal-derived nutrients in a porous microplate**. *The ISME Journal* 2022,

240        **16**(4):1036-1045.

241    17.    Lozupone CA, Hamady M, Kelley ST, Knight R: **Quantitative and qualitative beta diversity**

242        **measures lead to different insights into factors that structure microbial communities**. *Appl*

243        *Environ Microbiol* 2007, **73**(5):1576-1585.

244