

Nonparametric MANOVA via Independence Testing

Sambit Panda

Cencheng Shen, Ronan Perry, Jelle Zorn, Antoine Lutz, Carey E. Priebe,
Joshua T. Vogelstein

19 January, 2022

Johns Hopkins University

Department of Biomedical Engineering



Motivation

- Understand the relationship between k groups (*i.e.* control vs. disease)
- Question: Are they related? How?

One often desires to test groups

<i>U</i>	<i>V</i>
my grass	neighbor's grass
human brain connectivity	alien brain connectivity
control	disease
cancer risk group 1	cancer risk group 2

One often desires to test groups

<i>U</i>	<i>V</i>
my grass	neighbor's grass
human brain connectivity	alien brain connectivity
control	disease
cancer risk group 1	cancer risk group 2
any group	any other group

Statistics Background

- X is a random variable (some measurement)
- F_X is the distribution of X
- This means $F_X(a) = P(X \leq a)$
- This is denoted:

$$X \sim F_X$$

Statistics Background

- For two random variables X and Y , F_{XY} is called the joint distribution
- This means $F_{XY}(a, b) = P(X \leq a \text{ and } Y \leq b)$

or,

$$(X, Y) \sim F_{XY}$$

Informal Definition of Hypothesis Testing

- **Null Hypothesis:** The conventional belief about a phenomenon of interest, written H_0 .
- **Alternative Hypothesis:** An alternate belief about the same phenomenon, written H_A .
- **p-value:** The probability (under the null) of measurements more extreme than what was observed.

Formal Definition of K -Sample Testing

$$U_i^j \sim F_j, \quad j \in 1, \dots, k, \quad i \in 1, \dots, n_j$$

$$H_0 : F_1 = F_2 = \dots = F_k$$

$$H_A : \exists j \neq j' \text{ s.t. } F_j \neq F_{j'}$$

Note: These ideas and notation generalize for multivariate X and Y .

Outline

1. [Intuition](#)
2. [Simulations](#)
3. [Multiway and Multilevel](#)
4. [Real Data](#)
5. [Conclusion](#)

Intuition

Intuitive Desiderata of Testing Procedure

- Performant under *any* distribution
 - low- and high-dimensional
 - Euclidean and structured data (eg, sequences, images, networks, shapes)
 - linear and nonlinear relationships
- Is computational efficient

Provides a tractable algorithm that addresses the motivating question:

Are they related?

Analysis of Variance (ANOVA)

$$MST = \frac{\sum_{i=1}^k (T_i^2 / n_i) - G^2 / n}{k - 1}$$

$$MSE = \frac{\sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - \sum_{i=1}^k (T_i^2 / n_i)}{n - k}$$

$$ANOVA = \frac{MST}{MSE}$$

Multivariate ANOVA (MANOVA)

$$\mathbf{W} = \sum_{i=1}^k \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_{i.})(\mathbf{x}_{ij} - \bar{\mathbf{x}}_{i.})^T$$

$$\mathbf{B} = \sum_{i=1}^k n_i (\bar{\mathbf{x}}_{i.} - \bar{\mathbf{x}}_{..})(\bar{\mathbf{x}}_{i.} - \bar{\mathbf{x}}_{..})^T$$

$$MANOVA = \sum_{i=1}^s \frac{\lambda_i}{1 + \lambda_i} = \text{tr}(\mathbf{B}(\mathbf{B} + \mathbf{W})^{-1})$$

MANOVA Assumptions

- Data is derived from a multivariate Gaussian distribution
- Each group has the same covariance matrix

There must be a better test out there

Slight Tangent - Independence Testing

- X and Y are independent if neither contains information about the other
- In other words,

$$F_{XY} = P(X \leq a \text{ and } Y \leq b) = P(X \leq a) \times P(Y \leq b) = F_X F_Y$$

$$F_{XY} = F_X F_Y$$

Note: These ideas and notation generalize for multivariate X and Y .

Distance Correlation (Dcorr)

$$\widehat{Dcov}_{xy} = \frac{1}{n^2} \text{tr}(\mathbf{H}\mathbf{D}^x\mathbf{H}\mathbf{D}^y\mathbf{H})$$

$$\widehat{Dcorr}_{xy} = \frac{\widehat{Dcov}_{xy}}{\sqrt{\widehat{Dcov}_{xx} \times \widehat{Dcov}_{yy}}}$$

$$\mathbf{C}_{ij}^x = \mathbb{1}_{i \neq j} \left(\mathbf{D}_{ij}^x - \frac{1}{n-2} \sum_{t=1}^n \mathbf{D}_{it}^x - \frac{1}{n-2} \sum_{t=1}^n \mathbf{D}_{tj}^x + \frac{1}{(n-1)(n-2)} \sum_{t=1}^n \mathbf{D}_{tt}^x \right)$$

$$Dcov_{xy} = \frac{1}{n(n-3)} \text{tr}(\mathbf{C}^x\mathbf{C}^y)$$

$$Dcorr_{xy} = \frac{Dcov_{xy}}{\sqrt{Dcov_{xx} \times Dcov_{yy}}}$$

Multiscale Graph Correlation (MGC)

- Compute local Dcorr **at all scales**
- Find scale with **max** smoothed test statistic
- Permutation test to determine p-value

Kernel Mean Embedding Random Forest (KMERF)

- Train random forest on X , compute kernel matrix
- Transform similarity kernel matrix to distance matrix
- Permutation test to determine p-value

Great, what now?

- Can reduce the k -sample testing problem to the independence problem

$$\mathbf{x} = \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \\ \vdots \\ \mathbf{u}_k \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} \mathbf{1}_{n_1 \times 1} & \mathbf{0}_{n_1 \times 1} & \cdots & \mathbf{0}_{n_1 \times 1} \\ \mathbf{0}_{n_2 \times 1} & \mathbf{1}_{n_2 \times 1} & \cdots & \mathbf{0}_{n_2 \times 1} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_{n_k \times 1} & \mathbf{0}_{n_k \times 1} & \cdots & \mathbf{1}_{n_k \times 1} \end{bmatrix}$$

- Run any independence test
- ***Note: This process does not add any additional computational complexity to the independence testing algorithm***

Simulations

Definitions

- **power** is the probability of rejecting the null when the alternative is true

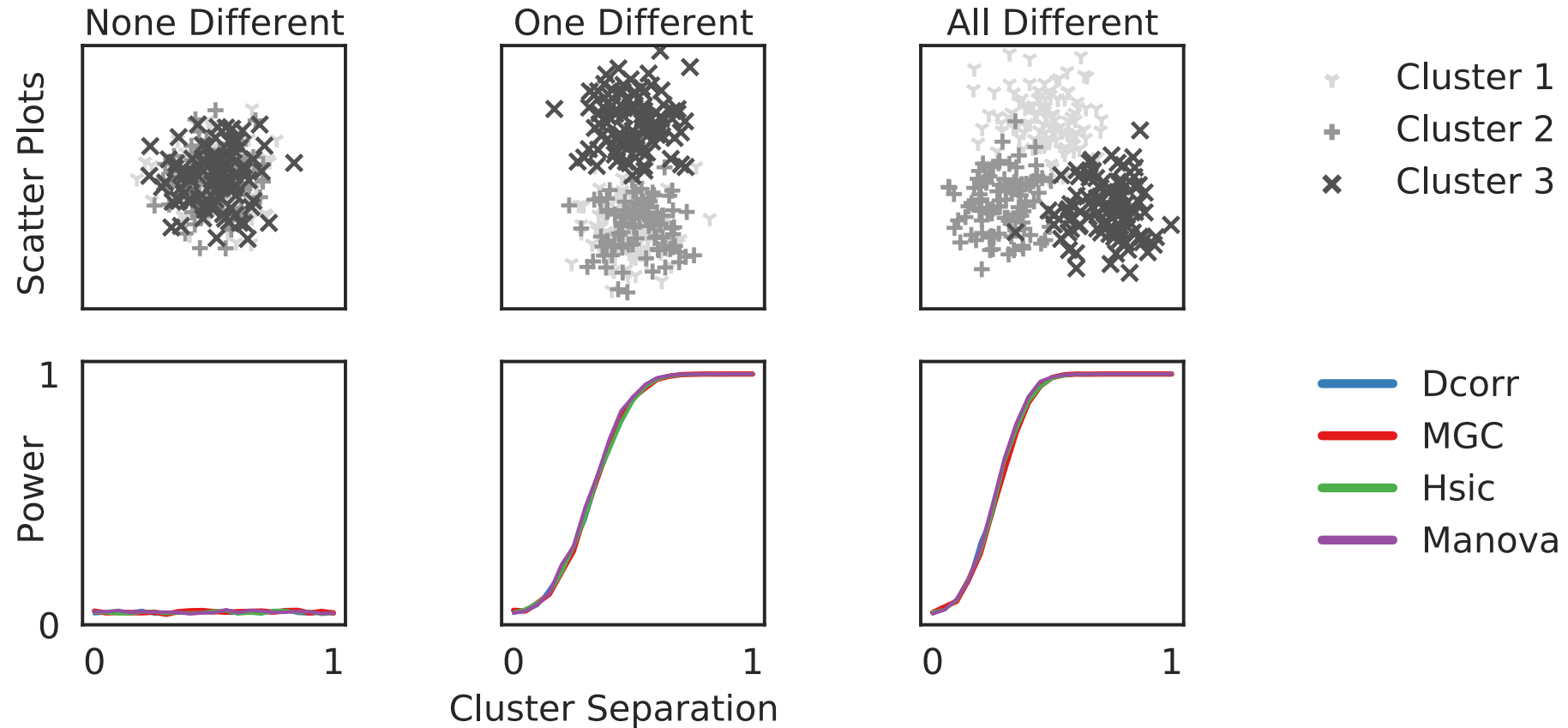
$\beta_n(t)$: power of test statistic t given n samples

- **relative power** power of one approach minus power of another

$$\beta_n(t) - \beta_n(\text{manova})$$

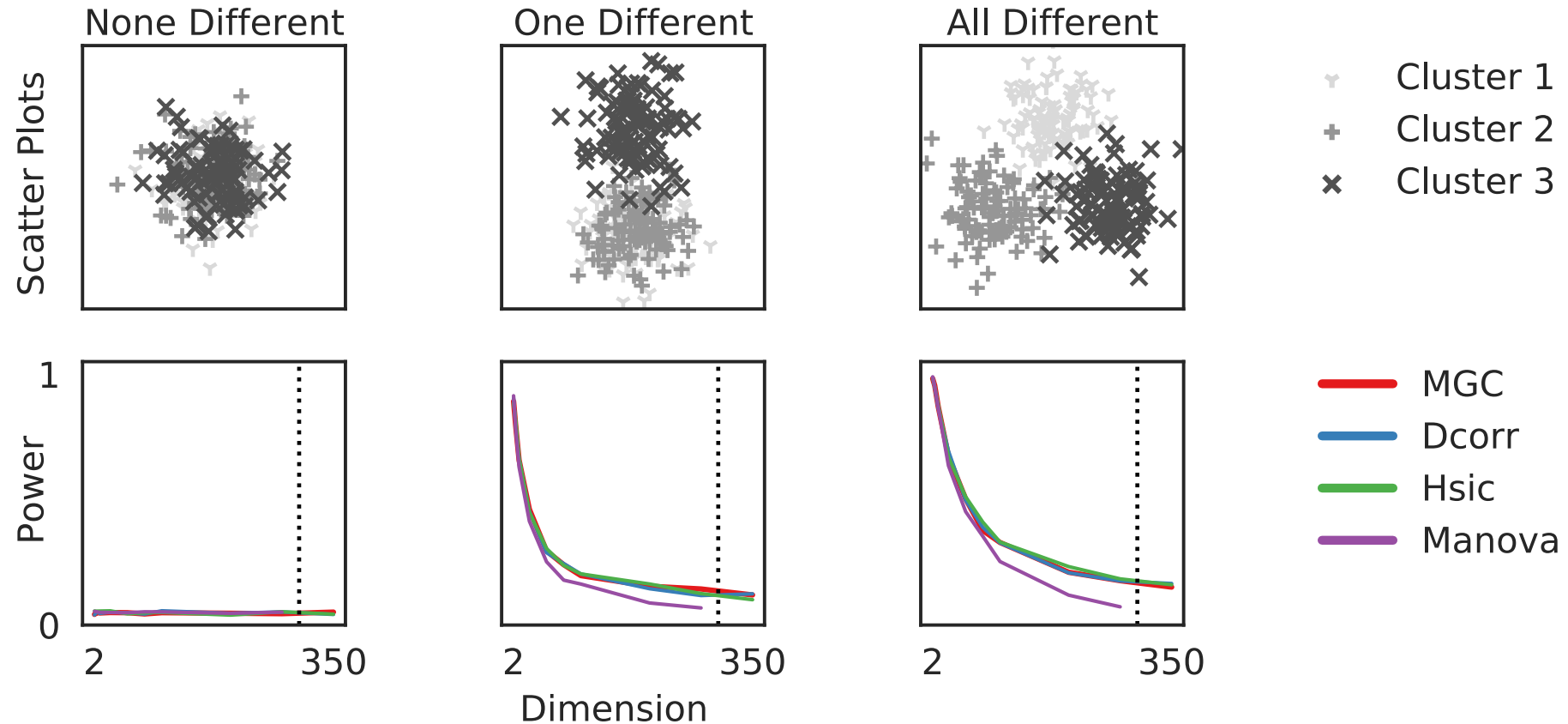
Optimal Settings for MANOVA (1D)

Power vs. increasing cluster separation

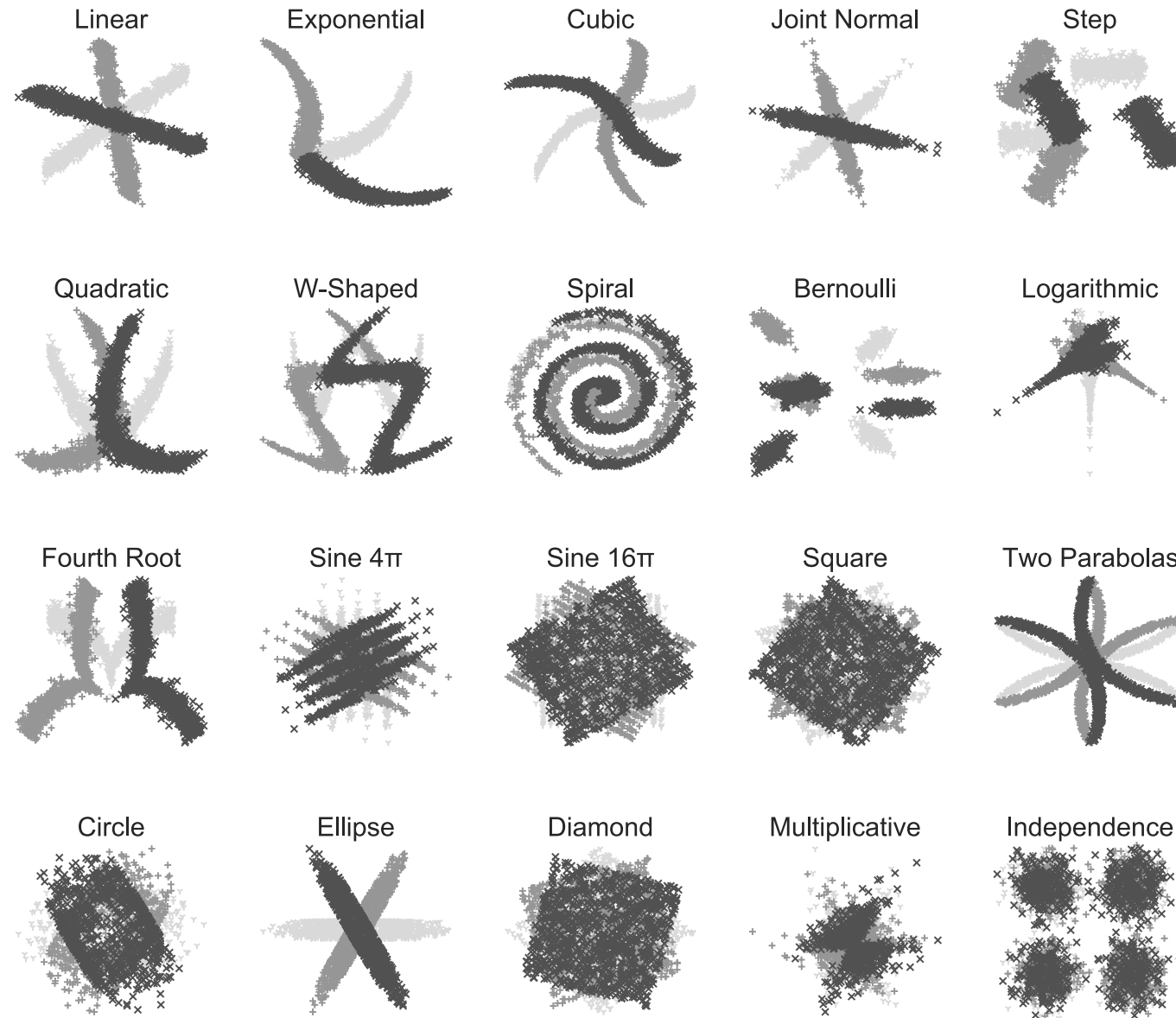


Optimal Settings for MANOVA (HD)

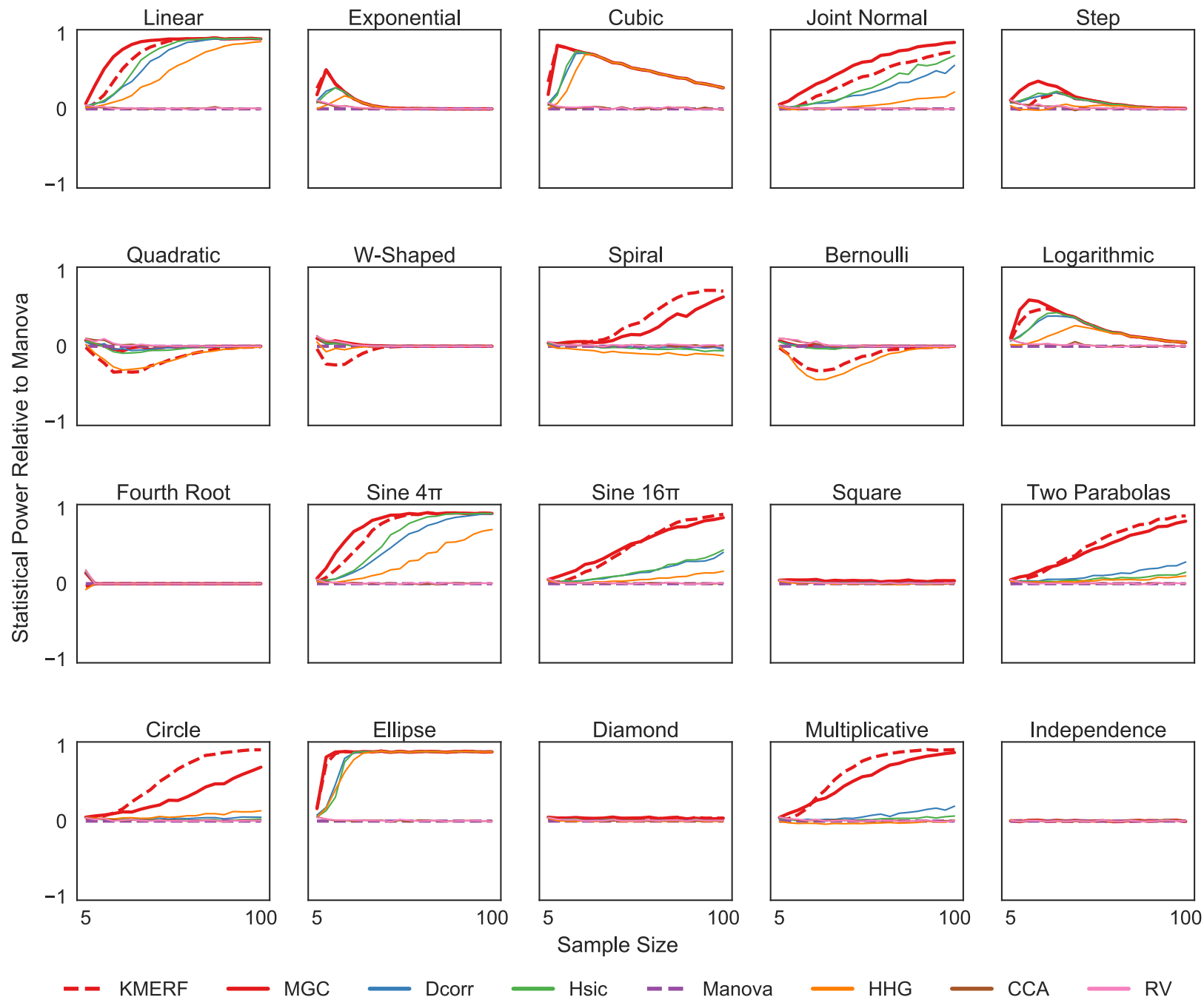
Power vs. increasing Gaussian dimension



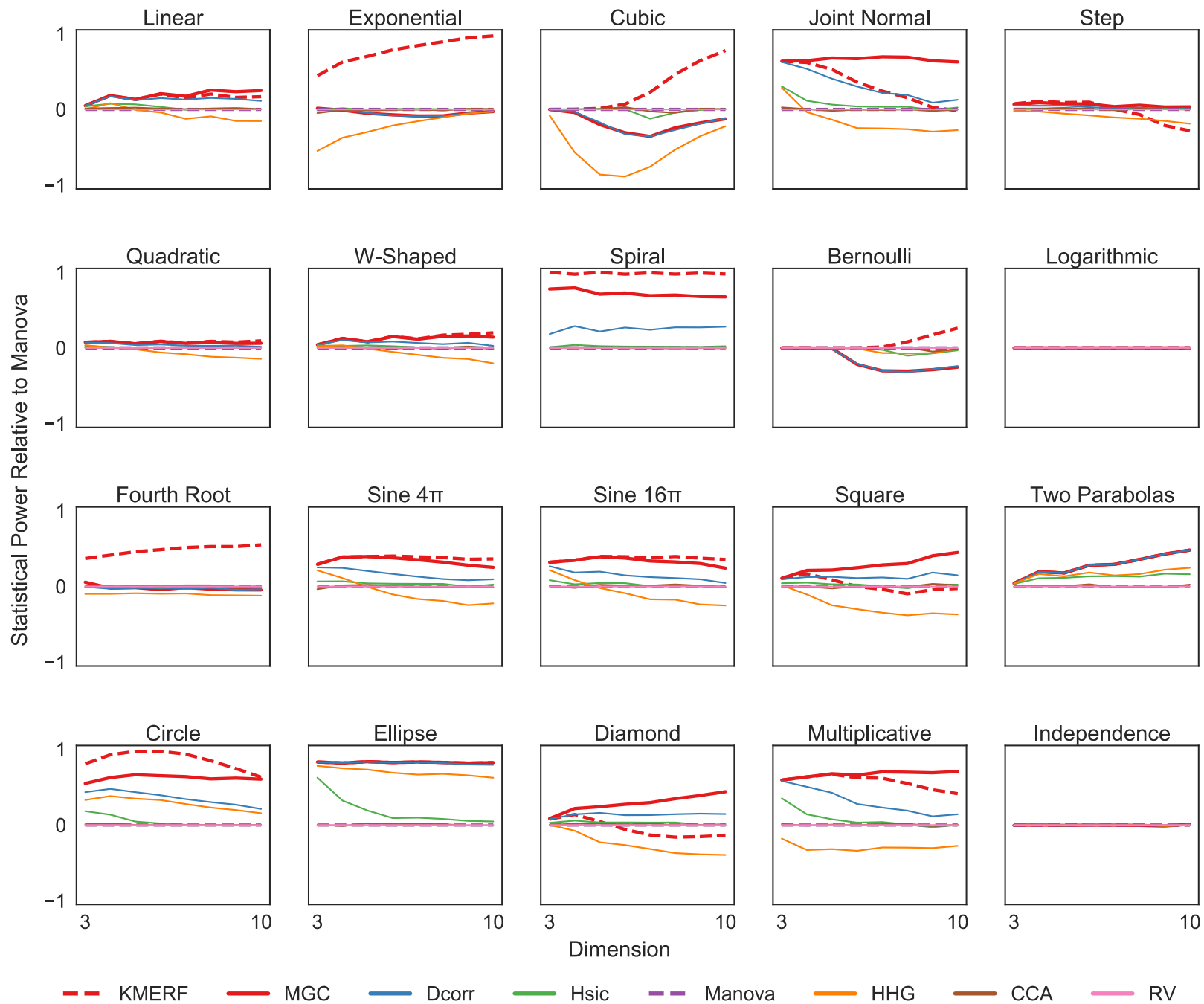
20 Different Functions (2D version)



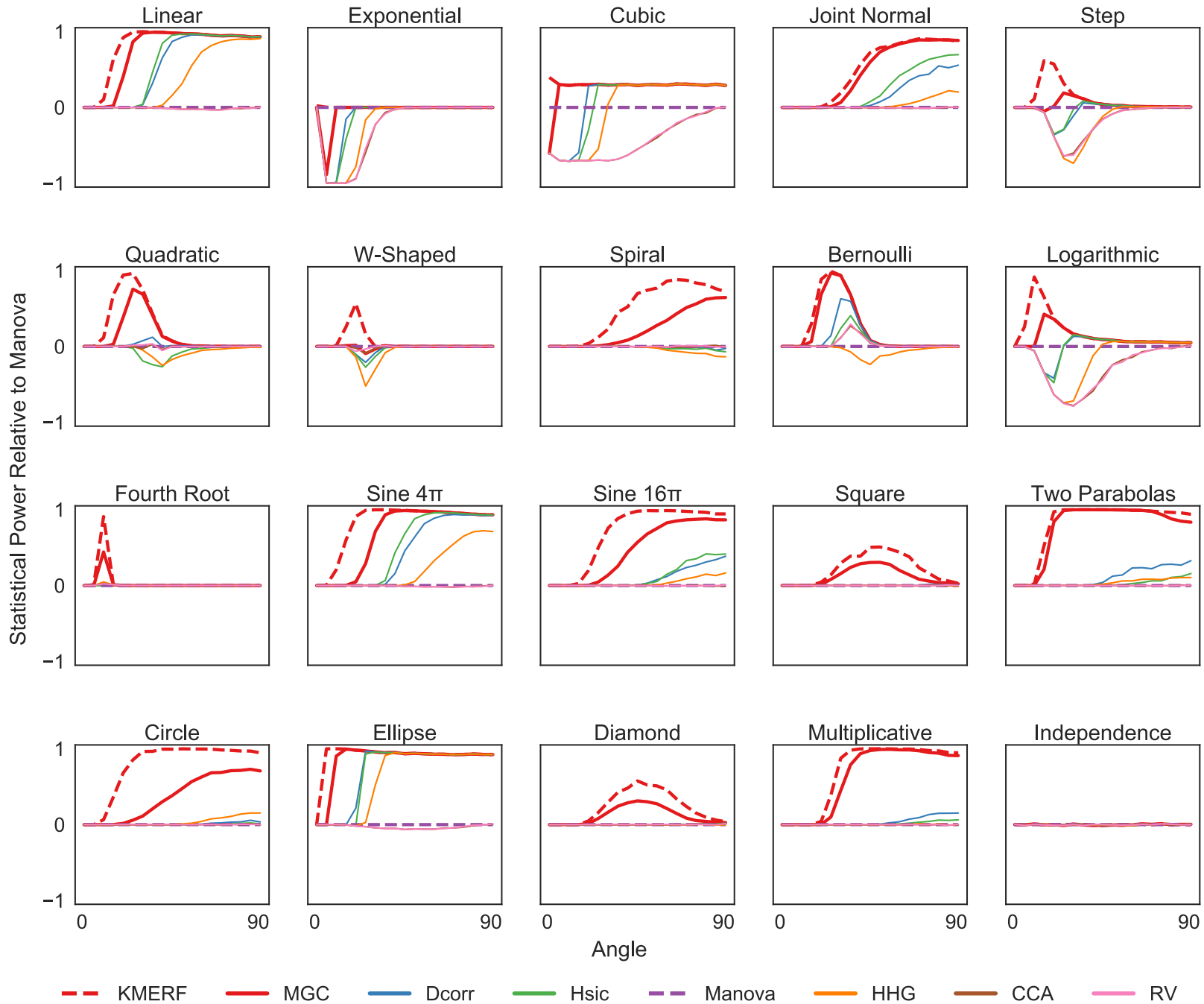
Multivariate Three-Sample Testing Increasing Sample Size



Multivariate Three-Sample Testing Increasing Dimension



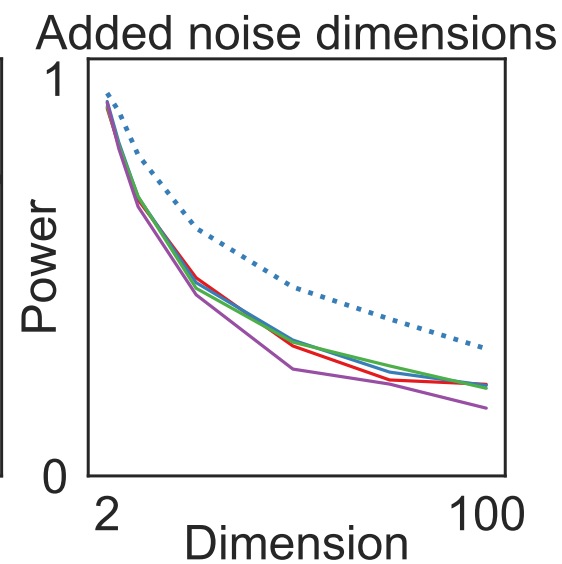
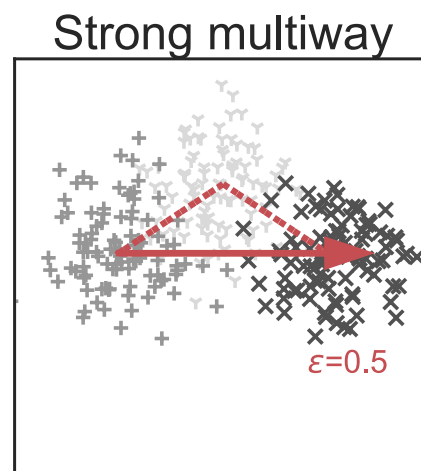
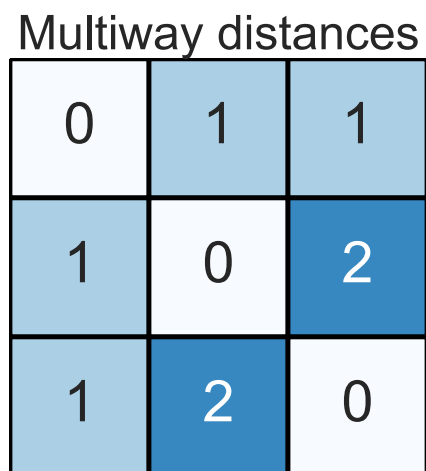
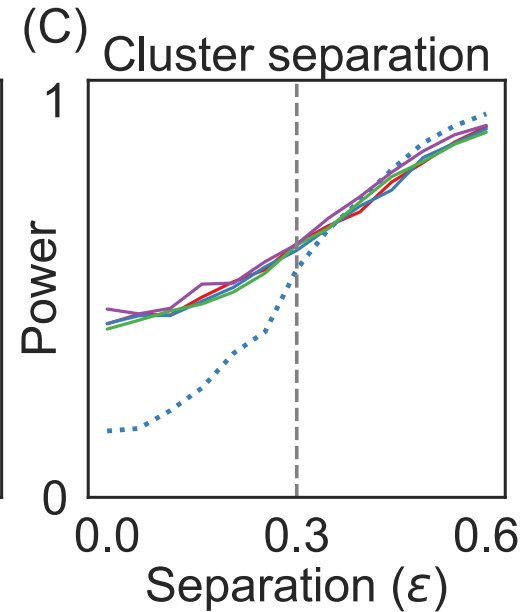
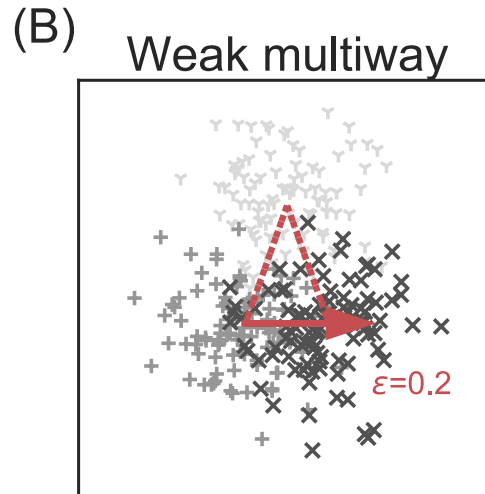
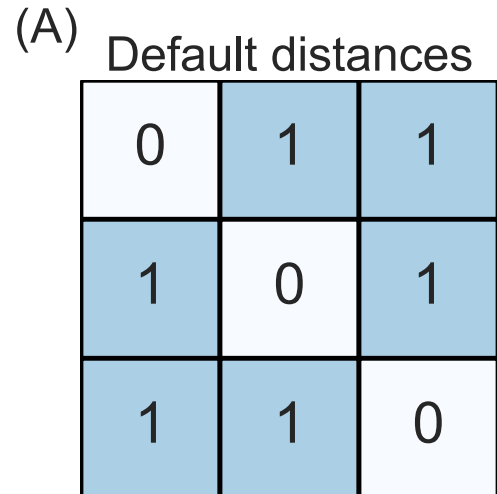
Multivariate Three-Sample Testing Increasing Angle



Multiway and Multilevel

Multiway Tests

- **Multiway:** More than one treatment group
- Instead of one-hot encoding, add 1's columns of label matrix

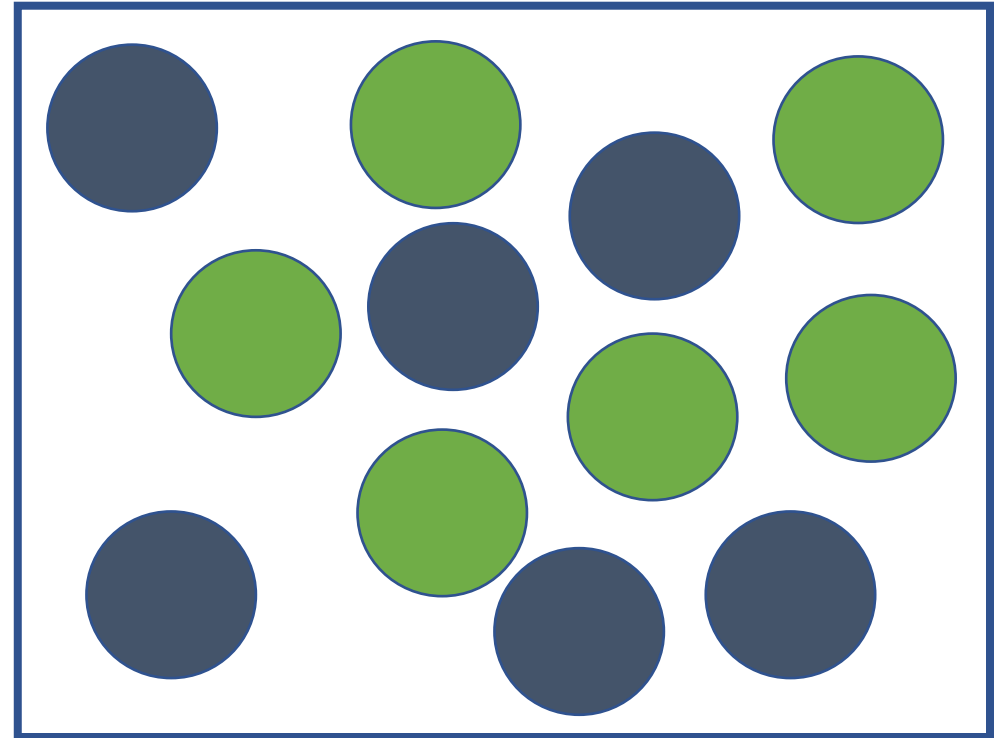
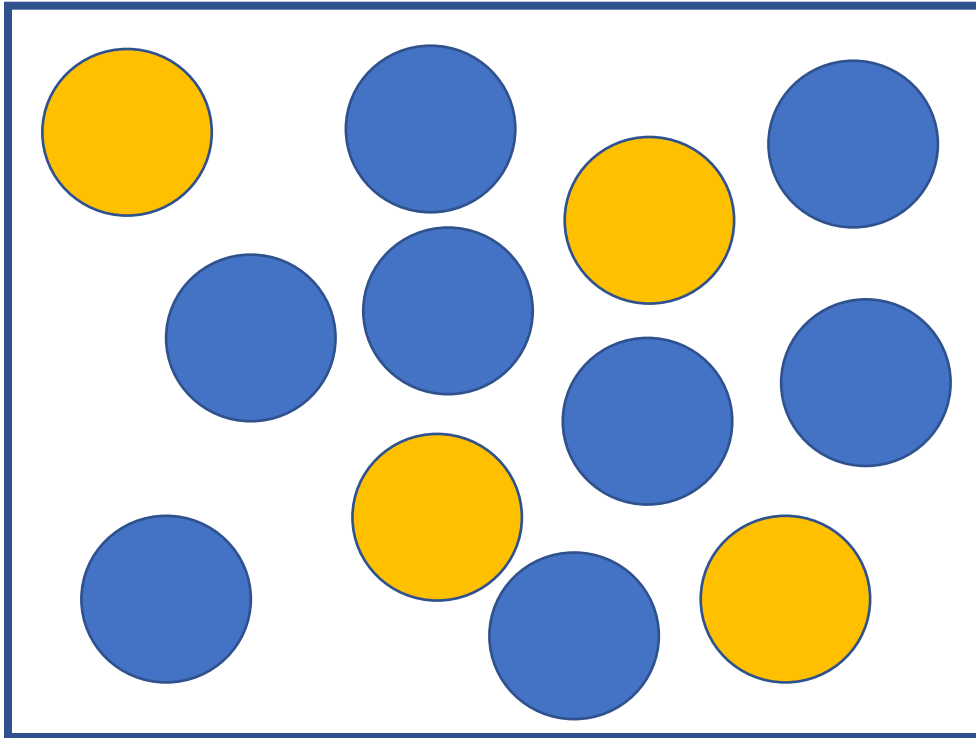


Multilevel Tests

- **Multilevel:** Samples are not always exchangeable with one another
- Need block permutation

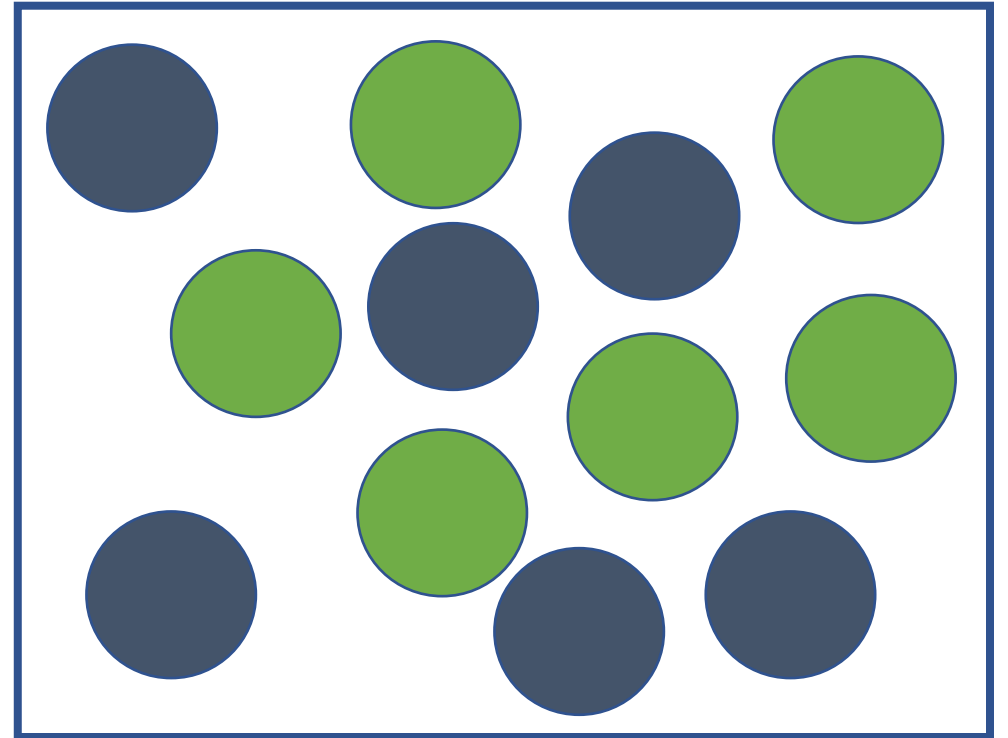
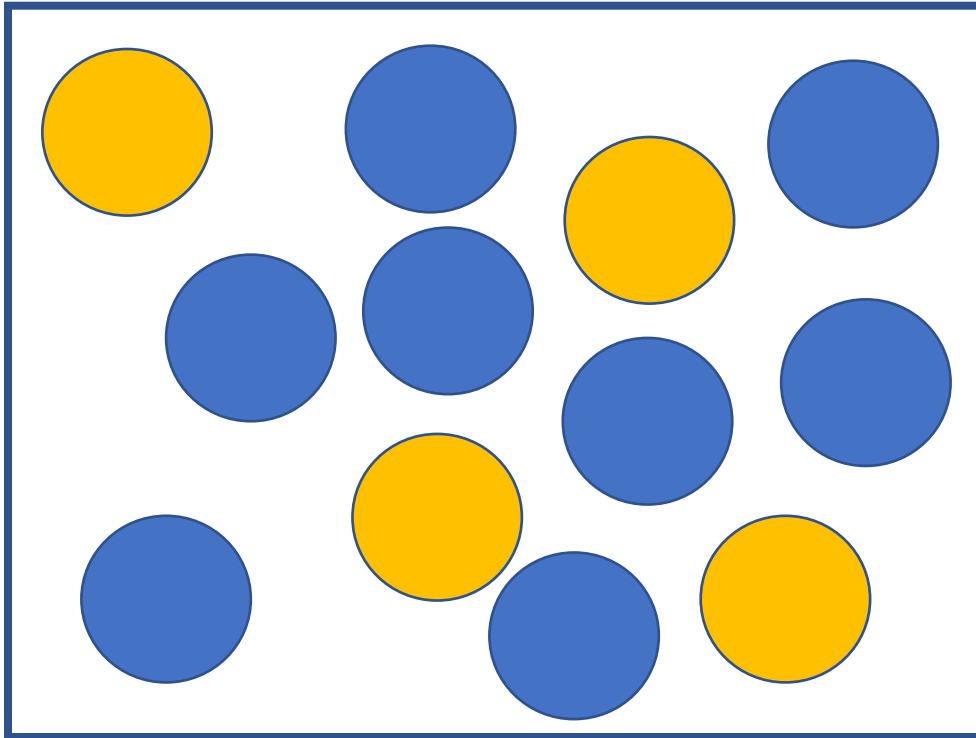
Multilevel Tests

- **Multilevel:** Samples are not always exchangeable with one another



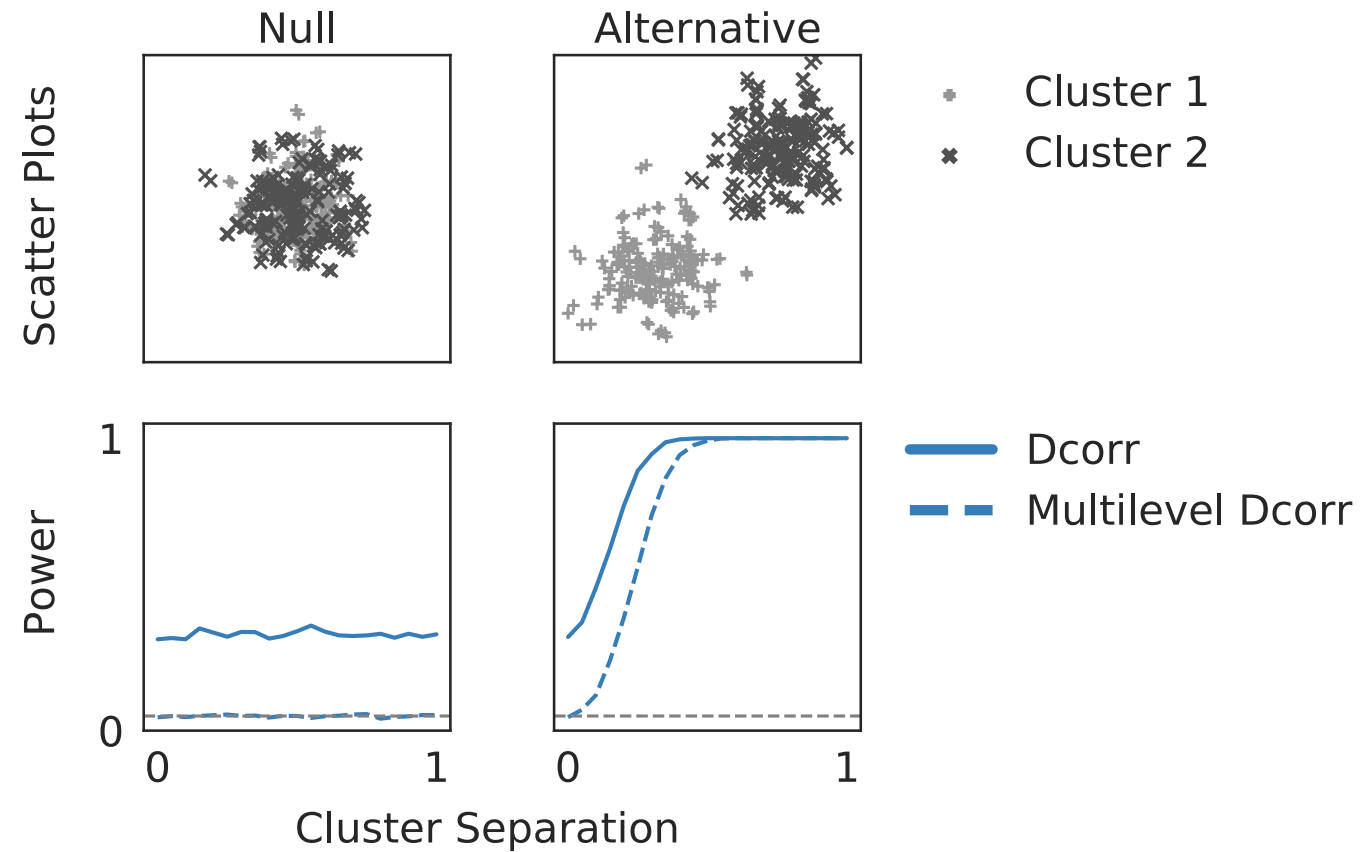
Multilevel Tests

- **Multilevel:** Samples are not always exchangeable with one another



Multilevel Tests

Multilevel Dcorr: Power vs. Cluster separation



Real Data

The Procedure

- Data: 75 subjects – 28 experienced and 47 novice meditators
- 3 Recording sessions for each meditator
- Computed gradients and tested for difference between traits and novice
- This is a **multilevel and multiway test**

Conclusions

- Presented several new k -sample tests using our framework
- At a simulation setting that fulfills MANOVA assumptions, our implementation performs as well or better
- Multiway tests give additional power when strong multiway effect is suspected
- Multilevel tests can now be performed

Next Steps

- All algorithms can be found in the [hyppo](#) package
 - [Documentation](#)
 - [Install](#)
 - [Tutorials](#)
- [Paper](#)

[Email](#) | [Website](#) | [Twitter](#)

Acknowledgements

- **Joshua Vogelstein, Cencheng Shen:** Theory, and paper writing
- **Ronan Perry:** Multiway, Multilevel, and Real Data
- **Jelle Zorn, Antoine Lutz:** Raw real data
- **Carey E. Priebe:** Theory
- **Russell Lyons, Minh Tang, Ronak Mehta, Eric Bridgeford:** Review
- ...and the rest of the NeuroData Lab



Questions?