

## Supporting Information 4: Eliminating accidental deviations to minimize generalization error and maximize replicability: applications in connectomics and genomics

Eric W. Bridgeford<sup>1</sup>, Shangsi Wang<sup>1</sup>, Zeyi Wang<sup>1</sup>, Ting Xu<sup>3</sup>, Cameron Craddock<sup>3</sup>, Jayanta Dey<sup>1</sup>, Gregory Kiar<sup>1</sup>, William Gray-Roncal<sup>1</sup>, Carlo Colantuoni<sup>1</sup>, Christopher Douville<sup>1</sup>, Stephanie Noble<sup>4</sup>, Carey E. Priebe<sup>1</sup>, Brian Caffo<sup>1</sup>, Michael Milham<sup>3</sup>, Xi-Nian Zuo<sup>2,5</sup>, Consortium for Reliability and Reproducibility, Joshua T. Vogelstein<sup>1,6\*</sup>

**S4 Simulations** The following simulations were constructed, where  $\sigma_{min}, \sigma_{max}$  are the variance ranges, and settings were run at 15 intervals in  $[\sigma_{min}, \sigma_{max}]$  for 500 repetitions per setting. For a simulation setting with variance  $\sigma$ , the variance is reported as the normalized variance,  $\bar{\sigma} = \frac{\sigma - \sigma_{min}}{\sigma_{max} - \sigma_{min}}$ . Dimensionality is 2, the number of items is  $K$ , and the total number of measurements across all items is 128. Typically,  $i$  indicates the individual identifier, and  $j$  the measurement index. Notationally, in the below descriptions, we adopt the convention that  $\mathbf{z}_i^j$  obeys the true distribution for a single observation  $j$  of item  $i$ , and  $\mathbf{x}_i^j$  incorporates the controlled error term  $\epsilon_i^j$ , which is the term which is varied the simulation. Further, each item features  $\frac{n}{K}$  measurements.

### Goodness of Fit Testing and Bayes Error

1. No Signal:  $K = 2$  items, where the true distributions for class 1 and class 2 are the same.
  - $\mathbf{z}_i^j \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}), i = 1, \dots, 2, t = 1, \dots, 64$ . Note:  $\mathbf{0} \in \mathbb{R}^2$  is  $\mathbf{0}$ , and likewise for  $\mathbf{I}$
  - $\epsilon_i^j \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}), \sigma \in [0, 20]$
  - $\mathbf{x}_i^j = \mathbf{z}_i^j + \epsilon_i^j \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, (1 + \sigma^2) \mathbf{I})$
2. Cross:  $K = 2$  items, where the true distributions for class 1 and class 2 are orthogonal.
  - $\Sigma_1 = \begin{bmatrix} 2 & 0 \\ 0 & 0.1 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 0.1 & 0 \\ 0 & 2 \end{bmatrix}$
  - $\mathbf{z}_i^j \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \Sigma_i), i = 1, 2$
  - $\epsilon_i^j \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}), \sigma \in [0, 20]$
  - $\mathbf{x}_i^j = \mathbf{z}_i^j + \epsilon_i^j$
3. Gaussian:  $K = 16$  items, where the true distributions are each gaussian.
  - $\mu_i \stackrel{iid}{\sim} \pi_1 \mathcal{N}(\mathbf{0}, 4\mathbf{I}), i = 1, \dots, 16$
  - $\Sigma = \begin{bmatrix} 1 & 0.1 \\ 0.1 & 1 \end{bmatrix}$
  - $\mathbf{z}_i^j \stackrel{iid}{\sim} \mathcal{N}(\mu_i, \Sigma)$
  - $\epsilon_i^j \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}), \sigma \in [0, 20]$
  - $\mathbf{x}_i^j = \mathbf{z}_i^j + \epsilon_i^j$
4. Ball/Circle:  $K = 2$  items, where 1 item is uniformly distributed on the unit ball with gaussian error, and the second item is uniformly distributed on the unit sphere with gaussian error.
  - $\mathbf{z}_1^t \stackrel{iid}{\sim} \mathbb{B}(r = 1) + \mathcal{N}(\mathbf{0}, 0.1\mathbf{I})$  samples uniformly on unit ball of radius 2 with Gaussian error
  - $\mathbf{z}_2^t \stackrel{iid}{\sim} \mathbb{S}(r = 1.5) + \mathcal{N}(\mathbf{0}, 0.1\mathbf{I})$  samples uniformly on unit sphere of radius 2 with Gaussian error

<sup>1</sup> Johns Hopkins University, Baltimore, Maryland, USA, <sup>2</sup> Shanghai Jiaotong University, Shanghai, China <sup>3</sup> Child Mind Institute, New York, New York, USA <sup>4</sup> Yale University, New Haven, Connecticut, USA <sup>5</sup> Beijing Normal University, Beijing, China, Nanning Normal University, Nanning, China, University of Chinese Academy of Sciences, Beijing, China, <sup>6</sup> Progressive Learning, Baltimore, Maryland, USA. \* [jovo@jhu.edu](mailto:jovo@jhu.edu).

- $\epsilon_i^j \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}), \sigma \in [0, 10]$
- $\mathbf{x}_i^j = \mathbf{z}_i^j + \epsilon_i^j$

5. XOR:  $K = 2$  items, where:

- $\mathbf{z}_1^t = \begin{cases} \mathbf{0} & t \in 1, \dots, 32 \\ \mathbf{1} & t \in 33, \dots, 64 \end{cases}$
- $\mathbf{z}_2^t = \begin{cases} [0, 1]' & t \in 1, \dots, 32 \\ [1, 0]' & t \in 33, \dots, 64 \end{cases}$
- $\epsilon_i^j \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}), \sigma \in [0, 0.8]$
- $\mathbf{x}_i^j = \mathbf{z}_i^j + \epsilon_i^j$

Bayes error was estimated by simulating  $n = 10,000$  points according to the above simulation settings, and approximating the Bayes error through numerical integration. The classification labels for  $K = 2$  simulations were consistent with the individual labels, and for the  $K = 16$ , the first class consists of the 8 distributions whose means were leftmost, and the rest of the distributions were the other class.

**Comparison Testing** Items are sampled with the same true distributions  $\mathbf{z}_i^j$  as before, with the following augmentation:

$$\mathbf{x}_{i,k}^j = \begin{cases} \mathbf{z}_i^j & k = 1 \\ \mathbf{z}_i^j + \epsilon_i^j & k = 2 \end{cases}$$

That is, the observed data  $\mathbf{x}_{i,k}^j$  for item  $i$ , observation  $j$ , and sample  $k \in [2]$  is such that the first sample is distributed according to the true item distribution, and the second sample is distributed according to the true item distribution with an added noise term, where  $\epsilon_i^j \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ :

1. No Signal:  $K = 2$   
 $\sigma \in [0, 10]$
2. Cross:  $K = 2$   
 $\sigma \in [0, 1]$
3. Gaussian:  $K = 16$   
 $\sigma \in [0, 1]$
4. Ball/Circle:  $K = 2$   
 $\sigma \in [0, 1]$
5. XOR:  $K = 2$   
 $\mathbf{x}_{i,k}^j = \begin{cases} \mathbf{z}_i^j + \tau_i^j & k = 1 \\ \mathbf{z}_i^j + \tau_i^j + \epsilon_i^j & k = 2 \end{cases}$  where  $\tau_i^j \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, 0.1 \mathbf{I})$   
 $\sigma \in [0, 0.2]$

By construction, one would anticipate  $\text{Discr}$  of the first sample to exceed that of the second sample, as the second sample has additional error. Therefore, the natural hypothesis is:

$$H_0 : D^{(1)} = D^{(2)}, \quad H_A : D^{(1)} > D^{(2)}$$