

The Approximate Fisher Information Matrix for Multinomial Mixture Models

Andrew M. Raim and Nagaraj K. Neerchal
Department of Mathematics and Statistics,
University of Maryland, Baltimore County



Summary

- Mixture distributions are useful in many problems, but can also be difficult to work with
- Minglei Liu (2005, PhD Thesis) studied estimation in multinomial mixture models. Related work has been done by Morel & Neerchal (1993, 1998, 2005)
- Here we present one of the key results, a large cluster approximation to the FIM
- The FIM approximation for the general multinomial mixture was shown to be useful in the Fisher Scoring algorithm. Here we consider its direct usage in inference

Mixture multinomial model

- Suppose we have s multinomial populations

$$f(\mathbf{x} | \mathbf{p}_1, m), \dots, f(\mathbf{x} | \mathbf{p}_s, m), \quad \mathbf{p}_\ell = (p_{\ell 1}, \dots, p_{\ell k})$$

- If population ℓ occurs with proportion π_ℓ , and we draw X from the mixed population

$$\mathbf{X} \sim f_\theta(\mathbf{x}) = \sum_{\ell=1}^s \pi_\ell f(\mathbf{x} | \mathbf{p}_\ell, m), \quad \mathbf{x} \in \mathcal{X}, \quad \theta = (\mathbf{p}_1, \dots, \mathbf{p}_s, \boldsymbol{\pi})$$

- Mixture distributions are a natural way to deal with mixed populations. A housing satisfaction survey from J. R. Wilson (1989) is an example featuring multinomials

Non-metropolitan area				Metropolitan area			
Neighborhood	US	S	VS	Neighborhood	US	S	VS
1	3	2	0	1	0	4	1
2	3	2	0	2	0	5	1
3	0	5	0	3	0	3	2
				...			
17	4	1	0	17	4	1	0
18	5	0	0				

- Some related problems (✓ means mixtures are often applied here)

Classification: Given samples $(\mathbf{x}_1^{(\ell)}, \dots, \mathbf{x}_{n_\ell}^{(\ell)})$ from each population $\ell = 1, \dots, s$, classify a new observation \mathbf{x}

Discriminant Analysis: Given samples $(\mathbf{x}_1^{(\ell)}, \dots, \mathbf{x}_{n_\ell}^{(\ell)})$, find a rule to best distinguish between the s groups

- Clustering: Given a sample $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ from the mixed population, try to determine which observations belong to the same groups. If s is not known the problem is harder

- Modeling overdispersion: In usual inference problems such as point estimation, confidence intervals, and hypothesis testing, model a mixed population with a mixture model of simpler distributions to capture the differences between groups

References

- M. Liu, *Estimation for Finite Mixture Multinomial Models*, PhD Thesis, University of Maryland, Baltimore County, Department of Mathematics and Statistics, 2005.
- J.G. Morel and N.K. Nagaraj, *A Finite Mixture Distribution for Modelling Multinomial Extra Variation*, *Biometrika* 80 (1993), pp. 363–371.
- N.K. Neerchal and J.G. Morel, *Large Cluster Results for Two Parametric Multinomial Extra Variation Models*, *Journal of the American Statistical Association* 93 (1998), pp. 1078–1087.
- N.K. Neerchal and J.G. Morel, *An improved method for the computation of maximum likelihood estimates for multinomial overdispersion models*, *Computational Statistics & Data Analysis* 49 (2005), pp. 33–43.

The computational resources used for this work were provided by the UMBC High Performance Computing Facility at the University of Maryland, Baltimore County (UMBC). See www.umbc.edu/hpcf for information on the facility and its uses.

FIM Approximation

- The Fisher Information Matrix (FIM, “outer product” form)

$$\mathcal{I}(\theta) := \mathbb{E} \left[\left(\frac{\partial}{\partial \theta} \log f_\theta(X) \right) \left(\frac{\partial}{\partial \theta} \log f_\theta(X) \right)^T \right]$$

is extremely useful for inference and model selection

- But a closed form expression cannot be obtained for most mixture models
- For the mixture of multinomials, the expectation can be computed exactly by summing over the sample space (which is finite but grows quickly with k and m)

$$\mathcal{I}(\theta) = n \sum_{\mathbf{x} \in \mathcal{X}} \left\{ \frac{\partial}{\partial \theta} \log f_\theta(\mathbf{x}) \right\} \left\{ \frac{\partial}{\partial \theta} \log f_\theta(\mathbf{x}) \right\}^T f_\theta(\mathbf{x})$$

- Liu and Morel & Nagaraj found an approximation for $\mathcal{I}(\theta)$

$$\tilde{\mathcal{I}}(\theta) := \begin{pmatrix} \pi_1 \mathbf{F}_1 & & 0 \\ & \ddots & \\ 0 & & \pi_s \mathbf{F}_s & \mathbf{F}_\pi \end{pmatrix} \quad (sk - 1) \times (sk - 1)$$

$$\mathbf{F}_\ell = m \left[\text{diag}(p_{\ell 1}^{-1}, \dots, p_{\ell, k-1}^{-1}) - p_{\ell k}^{-1} \mathbf{1}\mathbf{1}^T \right] \quad (k - 1) \times (k - 1)$$

$$\mathbf{F}_\pi = \text{diag}(\pi_1^{-1}, \dots, \pi_s^{-1}) - \pi_s^{-1} \mathbf{1}\mathbf{1}^T \quad (s - 1) \times (s - 1)$$

- This has a simple closed form that requires little computation to construct. There are also simple forms for the inverse approximate FIM and determinant
- Has been shown that $\tilde{\mathcal{I}}(\theta) - \mathcal{I}(\theta) \rightarrow 0$ as $m \rightarrow \infty$, for multinomial mixtures
- The approximation can also be used for more complicated mixtures such as Random-Clumped Multinomial and Dirichlet Multinomial. These multinomial mixtures feature parameters with functional dependencies on each other. See Neerchal & Morel (1998)
- In some cases the approximation was shown to work very well even for moderate m

Relationship to complete data FIM

- $\tilde{\mathcal{I}}(\theta)$ turns out to be equivalent to the complete data FIM obtained by considering latent class variables
- This technique is also used in Expectation Maximization (EM)
- Suppose we observe an iid sample $(\mathbf{X}_i, Z_i), i = 1, \dots, n$, where

$$Z_i = \begin{cases} 1 & \text{wp } \pi_1 \\ \dots & \\ s & \text{wp } \pi_s, \end{cases} \quad \mathbf{X}_i | Z_i = \ell \sim \text{Mult}(\mathbf{p}_\ell, m)$$

- Then the complete data likelihood is

$$f_\theta(\mathbf{x}, \mathbf{z}) = \prod_{i=1}^n \prod_{\ell=1}^s \left[\pi_\ell f(\mathbf{x}_i | \mathbf{p}_\ell, m) \right]^{I(z_i=\ell)}$$

- Computing the FIM (“Hessian” form) with respect to this likelihood, we obtain that

$$\mathbb{E} \left[- \frac{\partial^2}{\partial \theta \partial \theta^T} \log f_\theta(\mathbf{x}, \mathbf{z}) \right] \equiv \tilde{\mathcal{I}}(\theta)$$

- The classes (Z_1, \dots, Z_n) aren’t observable, but they are only used inside the expectation hence we don’t need to observe them

Simulation Study

The Wald statistic for testing $H_0 : \theta = \theta_0$ is

$$T_n(\hat{\theta}) = (\hat{\theta} - \theta_0)^T \mathcal{I}(\hat{\theta})(\hat{\theta} - \theta_0), \quad T_n(\hat{\theta}) \stackrel{L}{\rightarrow} \chi_q^2, \quad q = sk - 1.$$

This T can be used to construct an approximate $1 - \alpha$ level Wald-type confidence region

$$R(\hat{\theta}) = \left\{ \theta_0 : (\hat{\theta} - \theta_0)^T \mathcal{I}(\hat{\theta})(\hat{\theta} - \theta_0) \leq \chi_{q, \alpha}^2 \right\},$$

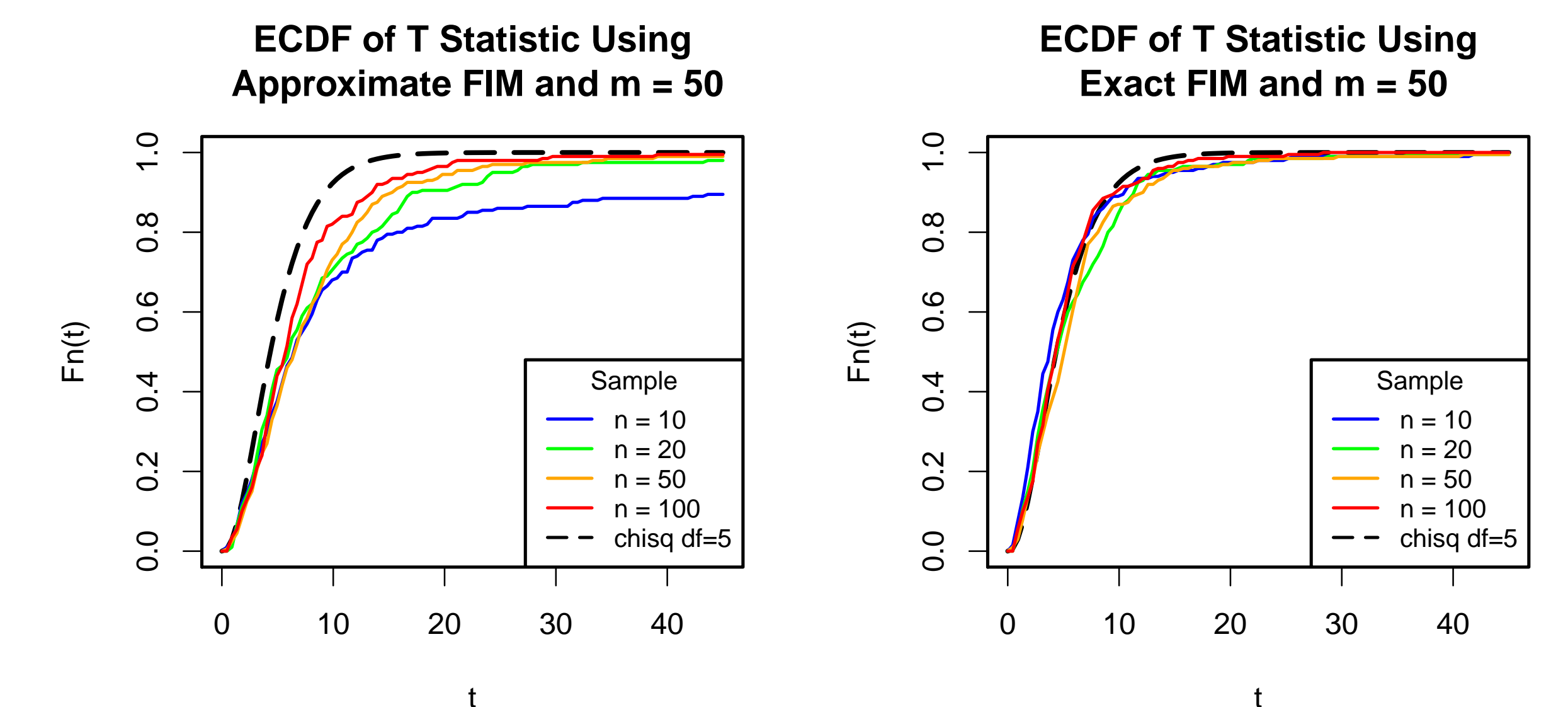
an ellipsoid in \mathbb{R}^q centered at the MLE $\hat{\theta}$, with shape determined by the FIM. Consider replacing $\mathcal{I}(\hat{\theta})$ with the approximate FIM, which is much easier to compute

$$\tilde{T}_n(\hat{\theta}) = (\hat{\theta} - \theta_0)^T \tilde{\mathcal{I}}(\hat{\theta})(\hat{\theta} - \theta_0).$$

We compare T and \tilde{T} with a simulation, choosing the parameters θ as

$$(\mathbf{p}_1 \quad \mathbf{p}_2 \quad \mathbf{p}_3) \propto \begin{pmatrix} 1 & 1 & 2 \\ 6 & 2 & 1 \end{pmatrix}, \quad \boldsymbol{\pi} \propto \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}.$$

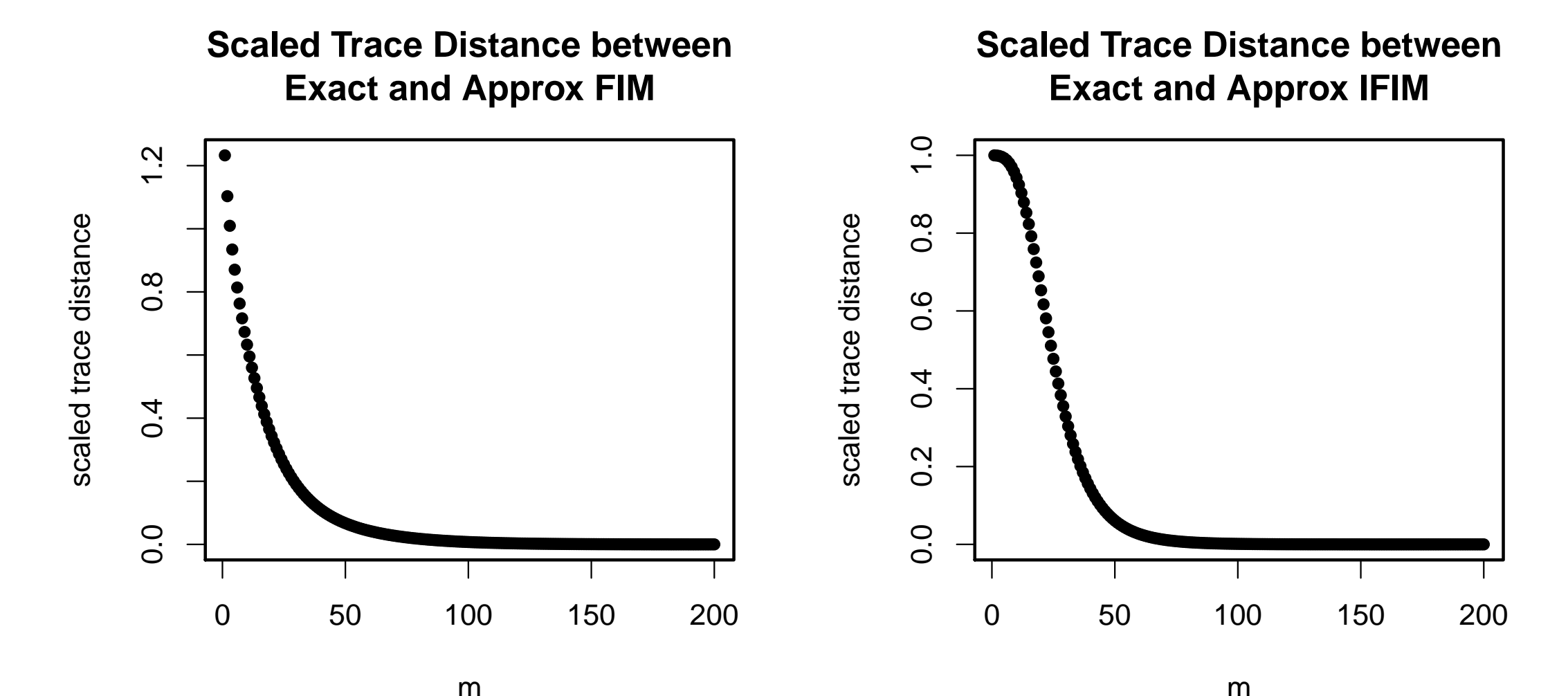
Samples were drawn from this binomial mixture 200 times for several m and n . For each sample we compute T and \tilde{T} , obtaining their empirical distributions under H_0 .



\tilde{T} seems to be lagging behind T in terms of the large sample χ_q^2 , even for $m = 50$ which may be considered a fairly large cluster size. To see why this is happening, we compare the two FIMs directly. Consider the following criteria based on the trace distance

$$d(A, B) = \frac{\text{tr}(A - B)^T(A - B)}{\text{tr } B^T B} = \frac{\sum_i \sum_j (a_{ij} - b_{ij})^2}{\sum_i \sum_j b_{ij}^2}.$$

For θ given above, we compute distances for varying m . The left plot shows $d(\tilde{\mathcal{I}}(\theta), \mathcal{I}(\theta))$ and the right plot shows $d(\tilde{\mathcal{I}}^{-1}(\theta), \mathcal{I}^{-1}(\theta))$ which corresponds to the asymptotic covariance matrix.



- Inference based on $\tilde{\mathcal{I}}(\theta)$ may not be correct for small to moderate m in the general mixture considered here
- Interesting that $\tilde{\mathcal{I}}(\theta)$ works well in estimation procedures like Fisher Scoring, even when m is not large. But it may be too far from $\mathcal{I}(\theta)$ to work well in inference
- $\tilde{\mathcal{I}}(\theta)$ may continue to be useful as a computational aid. Consider the approach of Neerchal and Morel (2005), where the approximation is used in Fisher Scoring iterations until convergence, and then one additional iteration is performed with the exact FIM to produce a final result