# Exploring Uniform Finite Sample Stickiness

Susanne Ulmer$^*$, Do Tran Van$^*$     Stephan F. Huckemann$^*$

May 18, 2023

### Abstract

It is well known, that Fréchet means on non-Euclidean spaces may exhibit nonstandard asymptotic rates depending on curvature. Even for distributions featuring standard asymptotic rates, there are non-Euclidean effects, altering finite sampling rates up to considerable sample sizes. These effects can be measured by the variance modulation function proposed by Pennec (2019). Among others, in view of statistical inference, it is important to bound this function on intervals of sampling sizes. In a first step into this direction, for the special case of a K-spider we give such an interval, based only on folded moments and total probabilities of spider legs and illustrate the method by simulations.

## 1 Introduction to Stickiness and Finite Sample Stickiness

Data analysis has become an integral part of science due to the growing amount of data in almost every research field. This includes a plethora of data objects that do not take values in Euclidean spaces, but rather in a non-manifold, stratified space. For statistical analysis in such spaces, it is therefore necessary to develop probabilistic concepts. Fréchet (1948) was one of the first to generalize the concept of an expected value to a random variable $X$ on an arbitrary metric space $(Q, \mathbf{d})$ as an minimizer of the expected squared distance:

$$\mu = \underset{p \in Q}{\operatorname{argmin}} \, \mathbb{E}[d(X, p)^2],  \tag{1}$$

nowadays called a *Fréchet mean* in his honor. Accordingly for a sample $X_1, \ldots, X_n \overset{\text{i.i.d.}}{\sim} X$, its Fréchet mean is given by

$$\mu_n = \underset{p \in Q}{\operatorname{argmin}} \sum_{j=1}^{n} \mathbf{d}^2(X_j, p).  \tag{2}$$

---

While on general space, these means can be empty or set valued, on *Hadamard spaces*, i.e. complete spaces of *global nonpositive curvature* (NPC), due to completeness, these means exists under very general conditions, and due to simple connectedness and nonpositive curvatures, they are unique, e.g. Sturm (2003), just as their Euclidean kin. They also share a law of strong numbers, i.e. that

$$\mu_n \overset{\text{a.s.}}{\to} \mu \,.$$

In contrast, however, their asymptotic distribution is often not normal, even worse, for some random variables their mean may be on a singular point stratum, and there may be a random sample size $N \in \mathbb{N}$ such that

$$\mu_n = \mu \text{ for all } n \geq N \,.$$

This phenomenon has been called *stickiness* by Hotz et al. (2013). It puts an end to statistical inference based on asymptotic fluctuation. While nonsticky means of random variables seem to feature the same asymptotic rate, as the Euclidean expected value, namely $1/\sqrt{n}$, it has been noted by Huckemann and Eltzner (2020) that for rather large sample sizes, the rates appear to be larger. This contribution is the first to systematically investigate this effect of *finite sample stickiness* on stratified spaces and we do this here for the model space of the $K$-spider introduced below. This effect is in some sense complementary to the effect of *finite sample smeariness*, where finite sample rates are smaller than $1/\sqrt{n}$. recently discovered by Hundrieser et al. (2020).

**Definition 1.** *With the above notation, assuming an existing Fréchet function $F(x) = \mathbb{E}[\mathbf{d}(X, x)^2] < \infty$ for all $x \in Q$, with existing Fréchet mean $\mu$,*

$$\mathbf{m}_n = \frac{n\mathbb{E}\left[\mathbf{d}^2(\mu_n, \mu)\right]}{\mathbb{E}\left[\mathbf{d}^2(X, \mu)\right]}, \tag{3}$$

*is the* variance modulation *for sample size n (see Pennec (2019)), or simply* modulation.

If $(Q, \mathbf{d})$ is Euclidean, then $\mathbf{m}_n = 1$ for all $n \in \mathbb{N}$, *smeariness* governs the cases $\mathbf{m}_n \to \infty$, see Hotz and Huckemann (2015); Eltzner and Huckemann (2019); Eltzner (2022), finite sample smeariness the cases $1 < \mathbf{m}_n$, cf. Tran et al. (2021); Eltzner et al. (2021, 2023), stickiness the case that $\mu_n = \mu$ a.s. for $n > N$ with a finite random sample size $N$, see Hotz et al. (2013); Huckemann et al. (2015); Barden et al. (2013, 2018), and *finite sample stickiness* the case that

$$0 < \mathbf{m}_n < 1 \text{ for nonsticky } \mu \,.$$

**Definition 2.** *For a nonsticky mean, if there are integers $l \in \mathbb{N}_{\geq 2}$, and $N \in N$ and $0 < \rho < 1$ such that*

$$0 < \mathbf{m}_n < 1 - \rho$$

*for all $n \in \{N, N + 1, ..., N^l\}$ then $X$ is called* finite sample sticky of level $\rho \in (0, 1)$, with scale $l$ and basis $N$ .
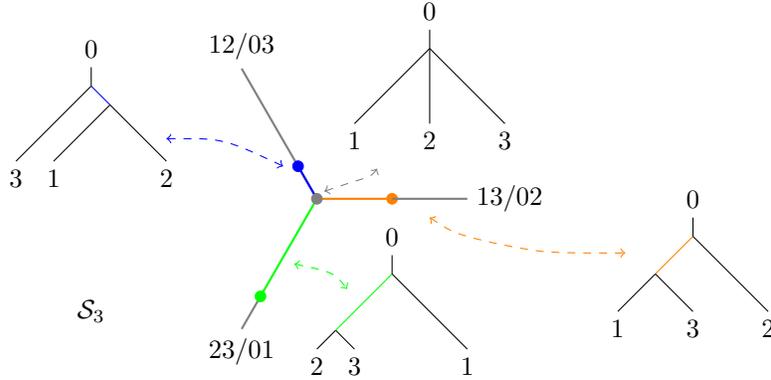
Figure 1: *Four different phylogenetic descendance trees for three taxa featuring one or none internal edge, modeled on the 3-Spider $\mathcal{S}_3$.*

We note that Pennec (2019) has shown that finite sample stickiness affects all affine connection manifolds with constant negative sectional curvature. We conjecture that this is also the case for general Hadamard spaces.

A very prominent example of a nonmanifold Hadamard space is given by the *BHV tree spaces* introduced by Billera et al. (2001) modeling phylogenetic descendance trees. For a fixed number of species or taxa the BHV-space models all different tree topologies, where within each topology, lengths of internal edges reflect evolutionary mutation from unknown ancestors. For three taxa, there are three topologies featuring nonzero internal edges and a fourth one, the *star tree* featuring no internal edge. The corresponding BHV space thus carries the structure of a 3-spider as depicted in Figure 1. For illustration of argument, in this contribution we consider $K$-spiders.

## 2   Model Space: The K-Spider

The following has been taken from Hotz et al. (2013).

**Definition 3.** *For $3 \leq K \in \mathbb{N}$ the $K$-spider $\mathcal{S}_K$ is the space*

$$\mathcal{S}_K = [0, \infty) \times \{1, 2, ..., K\} / \sim$$

*where for $i, k \in \{1, 2, ..., K\}$ and $x \geq 0$, $(x, i) \sim (x, k)$ if $k = i$ or $x = 0$. The equivalence class of $(0, 1)$ is identified with the* origin **0***, so that $\mathcal{S}_K = \{\mathbf{0}\} \cup \bigcup_{k=1}^{K} L_k$ with the positive half-line $L_k := (0, \infty) \times \{k\}$ called the $k$-th leg.*

*Further, for any $k \in \{1, 2, ..., K\}$ the map*

$$F_k : \mathcal{S}_K \to \mathbb{R},$$

$$(x, i) \mapsto \begin{cases} x & if \quad i = k, \\ -x & else \end{cases}$$

*is called the k-th folding map of $\mathcal{S}_K$.*

*The first two folded moments on the k-th leg of a random variable $X$ are*

$$m_k = \mathbb{E}[F_k(X)], \quad \sigma_k^2 = \mathbb{E}[(F_k(X) - m_k)^2]$$

*in population version and the first folded moment in sample version for random variables $X_1, \ldots, X_n$ is*

$$\eta_{n,k} = \frac{1}{n} \sum_{j=1}^{n} F_k(X_j) \,.$$

*We say the $X$ is* nondegnerate *if $\mathbb{P}\{X \in L_k\} > 0$ for at least three different $k \in \{1, \ldots, K\}$.*
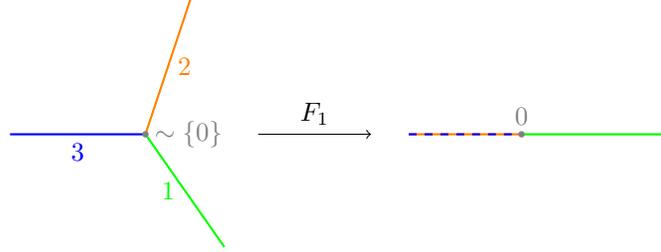


Figure 2: *The first folding map $F_1$ on the Three-Spider $\mathcal{S}_3$. The leg labelled 1 is mapped to the positive real line. The second and third leg are folded and then mapped to the negative half line.*

**Lemma 4.** *For a sample of size $n$ of a nondegenerate random variable, we have for every $1 \le k \le K$ that*

$$\eta_{n,k} > 0 \quad \Leftrightarrow \quad \mu_n \in L_k \quad \Leftrightarrow \quad F_k(\mu_n) > 0 \tag{4}$$

*whereas*

$$\eta_{n,k} = 0 \Rightarrow \mu_n = \mathbf{0} \in \mathcal{S}_K \Rightarrow \eta_{n,k} \le 0 \,.$$

*In particular:*

$$\eta_{n,k} \ge 0 \quad \Leftrightarrow \quad F_k(\mu_n) = \eta_{n,k} \tag{5}$$

*and*

$$\eta_{n,k} < 0 \quad \Leftrightarrow \quad \eta_{n,k} < F_k(\mu_n) \,. \tag{6}$$

4

*Proof.* All but the last assertion are from Hotz et al. (2013, Lemma 3.3). The last is also in the proof of Huckemann and Eltzner (2024, Lemma 4.1.4), we prove it here for convenience.

Letting $h_{n,i} = \frac{1}{n} \sum_{X_j \in L_i} F_i(X_j)$ for $i \in \{1, \ldots, K\}$ we have

$$\eta_{n,k} \;\; = \;\; h_{n,k} - \sum_{\substack{i \neq k \\ 1 \leq i \leq K}} h_{n,i} \, . \tag{7}$$

If $\eta_{n,k} < 0$ then there must be $k \neq k' \in \{1, \ldots, K\}$ with $\mu_n \in L_{k'}$ and thus

$$\begin{aligned}
0 \;\; < \;\; -F_k(\mu_n) \;\; = \;\; F_{k'}(\mu_n) \;\; = \;\; \eta_{n,k'} \;\; &= \;\; h_{n,k'} - h_{n,k} - \sum_{\substack{i \neq k, k' \\ 1 \leq i \leq K}} h_{n,i} \\
&= \;\; -\eta_{n,k} - 2 \sum_{\substack{i \neq k, k' \\ 1 \leq i \leq K}} h_{n,i} \;\; < \;\; -\eta_{n,k}
\end{aligned}$$

as asserted, due to nondegenaracy. $\qquad\qquad\square\qquad\qquad\qquad\square$

The case, $m_k < F_k(\mu)$ for all $k \in \{1, \ldots, K\}$ governs stickiness, cf. Hotz et al. (2013), and the case of $\eta_{n,k} < F_k(\mu_n) \leq 0 < F_k(\mu) = m_k$ for some $k \in \{1, \ldots, K\}$ governs finite sample stickiness, as detailed below. For this reason, we consider the following.

**Corollary 5.** *In case of $m_k > 0$, we have $|F_k(\mu_n) - m_k|^2 \leq |\eta_{n,k} - m_k|^2$ and $n\mathbb{E}[\mathbf{d}^2(\mu_n, \mu)] = n\mathbb{E}[|F_k(\mu_n) - m_k|^2] \leq \sigma_k^2$.*

*Proof.* If $\eta_{n,k} \geq 0$ we have by (5) that $F_k(\mu_n) = \eta_{n,k}$ so that $|F_k(\mu_n) - m_k|] = |F_k(\mu_n) - m_k|$. If $\eta_{n,k} < 0$ we have by (6) that $\eta_{n,k} < F_k(\mu_n)$ and hence $|\eta_{n,k} - m_k| = m_k - \eta_{n,k} > m_k - F_k(\mu_n) = |m_k - F_k(\mu_n)|$, where the last equality is due to (4).

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

# 3 Estimating Uniform Finite Sample Stickiness

Denote the standard-normal-cdf by $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{x} \exp(-t^2/2)dt$.

**Theorem 6.** *(Berry-Esseen, Esseen (1945, p. 42), Shevtsova (2011))*
*Let $Z_1, Z_2, \ldots, Z_n$ be iid random variables on $\mathbb{R}$ with mean zero and finite third moment $\mathbb{E}\left[|Z_1|^3\right] < \infty$. Denote by $\sigma^2$ the variance and by $\hat{F}_n(z)$ the cumulative distribution function of the random variable*

$$\frac{Z_1 + Z_2 + \ldots + Z_n}{\sqrt{n}\sigma}.$$

*Then there is a finite positive constant $C_S \leq 0.4748$ such that*

$$|\hat{F}_n(z) - \Phi(z)| \leq \frac{\mathbb{E}\left[|Z_1|^3\right]}{\sqrt{n}\sigma^3} C_S. \tag{8}$$

5

For all of the following let $X_1, \ldots, X_n \overset{\text{i.i.d.}}{\sim} X$ be nondegenerate random variables on $\mathcal{S}_K$ with Fréchet mean $\mu \in L_k$ for some $k \in \{1, \ldots, K\}$ and $\mathbb{E}\left[\mathbf{d}^3(X, \mathbf{0})\right] < \infty$. Let

$$p_n = \sum_{i=1}^{K} \Phi\left(\frac{\sqrt{n}m_i}{\sigma_i}\right) + \sum_{i=1}^{K} \frac{\mathbb{E}\left[|F_i(X) - m_i||^3\right]}{\sqrt{n}\sigma_i^3} C_S,$$

$$p_{n,k} = \Phi\left(\frac{\sqrt{n}m_k}{\sigma_k}\right) - \frac{\mathbb{E}\left[|F_k(X) - m_k||^3\right]}{\sqrt{n}\sigma_k^3} C_S.$$

For $i \in \{1, \ldots, K\}$ set

$$A_i = \{\eta_{n,i} = F_i(\mu_n)\}, \quad A = A_1 \cup \ldots \cup A_K.$$

Due nondegeneracy, (7) implies that the $A_i$ $(i = 1, \ldots, K)$ are disjoint, see also (Hotz et al., 2013, Theorem 2.9).

**Lemma 7.** *For all $n \in \mathbb{N}$, $\mathbb{P}(A_k) \geq p_{n,k}$ and $\mathbb{P}(A) \leq p_n$.*

*Proof.* Fix $i \in \{1, \ldots, K\}$. By Lemma 4, $\qquad\qquad\qquad\qquad\qquad\qquad\square$

$$\mathbb{P}(A_i) = \mathbb{P}\left\{\eta_{n,i} \geq 0\right\} = \mathbb{P}\left\{\frac{\sqrt{n}(-\eta_{n,i} + m_i)}{\sigma_i} \leq \frac{\sqrt{n}m_i}{\sigma_i}\right\}.$$

Setting $Z_j = F_i(X_j) - m_i$ $(j = 1, \ldots, n)$ and $z = \frac{\sqrt{n}m_i}{\sigma_i}$ we can apply the Berry-Esseen theorem, Theorem 6, yielding for $i = k$

$$\mathbb{P}(A_k) \geq \Phi\left(\frac{\sqrt{n}m_k}{\sigma_k}\right) - \frac{\mathbb{E}\left[|F_k(X) - m_k||^3\right]}{\sqrt{n}\sigma_k^3} C_S = p_{n,k}$$

and

$$\mathbb{P}(A) \leq \sum_{i=1}^{K} \Phi\left(\frac{\sqrt{n}m_i}{\sigma_i}\right) + \sum_{i=1}^{K} \frac{\mathbb{E}\left[|F_i(X) - m_i|^3\right]}{\sqrt{n}\sigma_i^3} C_S = p_n.$$

**Theorem 8.** *A nondegenerate random variable on $X$ on $\mathcal{S}_K$ with mean $\mu \in L_k$ and second folded moment $\sigma_k^2$, $k \in \{1, \ldots, K\}$, is finite sample sticky of level*

$$\rho = \min_{n \in \{N, N+1, \ldots, N^l\}} 1 - p_n - \frac{nm_k^2}{\sigma_k^2}(1 - p_{n,k})$$

*with scale $l$ and basis $N$, if there is $l \in \mathbb{N}_{\geq 2}$ such that $p_n < 1$ and $p_{n,k} \geq 0$, and if*

$$p_n + \frac{nm_k^2}{\sigma_k^2}(1 - p_{n,k}) < 1$$

*for all $n \in \{N, N+1, \ldots, N^l\}$.*

6

*Proof.* By the law of total expectation, Corollary 5 and Lemma 7, exploiting that $A^c \subseteq A_k^c$, we obtain

$$\mathbf{m}_n = \mathbb{E}[\mathbf{m}_n] = \frac{n}{\sigma_k^2} \mathbb{E}\left[\mathbb{E}[\mathbf{d}^2(\mu_n, \mu) \mid A]\mathbb{P}(A) + \mathbb{E}[\mathbf{d}^2(\mu_n, \mu) \mid A^c]\mathbb{P}(A^c)\right]$$

$$= \frac{n}{\sigma_k^2} \mathbb{E}[\mathbf{d}^2(\mu_n, \mu)]\mathbb{P}(A) + \frac{nm_k^2}{\sigma_k^2}\mathbb{P}(A^c)$$

$$\leq p_n + \frac{nm_k^2}{\sigma_k^2}(1 - p_{n,k}).$$

$\square$

# 4    Example and Simulations

**Example 9.** *For $t > 0$ let $X_t$ be a random variable on $\mathcal{S}_K$ with probabilities*

$$\mathbb{P}\left\{X_t = (K - 1 + Kt, K)\right\} = \frac{1}{K} = \mathbb{P}\left\{X_t = (1, i)\right\}, \qquad k = i, ..., K - 1\,,$$
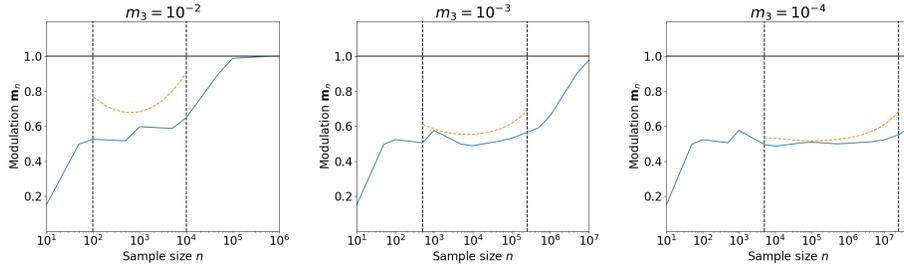
*and first moments*

$$m_i = \frac{4 - K(2 + t)}{K}, \qquad m_K = t.$$

*Hence $X_t$ is nonsticky with $\mu \in L_K$ for $t > 0$ and*

$$\mathbb{E}[\mathbf{d}^3(X_t, 0)] = \frac{K - 1}{K} + \frac{(K - 1 + Kt)^3}{K} < \infty.$$

*We can therefore apply Theorem 8. Figure 3 illustrates intervals of sample sizes displaying finite sample stickiness for the 3-Spider and $t \in \{10^{-2}, 10^{-3}, 10^{-4}\}$. The explicit bound for the modulation derived by Theorem 8 is given as an orange dashed line.*



(a) *Finite sample stickiness of level $\rho = 0.1$ with scale 2 and base $N = 100$.*

(b) *Finite sample stickiness of level $\rho = 0.31$ with scale 2 and base $N = 500$.*

(c) *Finite sample stickiness of level $\rho = 0.31$ with scale 2 and base $N = 5000$.*

Figure 3

# Acknowledgment

# References

Barden, D., H. Le, and M. Owen (2013). Central limit theorems for Fréchet means in the space of phylogenetic trees. *Electron. J. Probab 18*(25), 1–25.

Barden, D., H. Le, and M. Owen (2018). Limiting behaviour of Fréchet means in the space of phylogenetic trees. *Annals of the Institute of Statistical Mathematics 70*(1), 99–129.

Billera, L., S. Holmes, and K. Vogtmann (2001). Geometry of the space of phylogenetic trees. *Advances in Applied Mathematics 27*(4), 733–767.

Eltzner, B. (2022). Geometrical smeariness–a new phenomenon of Fréchet means. *Bernoulli 28*(1), 239–254.

Eltzner, B., P. Hansen, S. F. Huckemann, and S. Sommer (2023). Diffusion means in geometric spaces. *Bernoulli*. to appear.

Eltzner, B. and S. F. Huckemann (2019). A smeary central limit theorem for manifolds with application to high-dimensional spheres. *Ann. Statist. 47*(6), 3360–3381.

Eltzner, B., S. Hundrieser, and S. Huckemann (2021). Finite sample smeariness on spheres. In *Geometric Science of Information: 5th International Conference, GSI 2021, Paris, France, July 21–23, 2021, Proceedings 5*, pp. 12–19. Springer.

Esseen, C.-G. (1945). Fourier analysis of distribution functions. a mathematical study of the Laplace-Gaussian law. *Acta Mathematica 77*, 1–125.

Fréchet, M. (1948). Les éléments aléatoires de nature quelconque dans un espace distancié. *Annales de l'Institut de Henri Poincaré 10*(4), 215–310.

Hotz, T. and S. Huckemann (2015). Intrinsic means on the circle: Uniqueness, locus and asymptotics. *Annals of the Institute of Statistical Mathematics 67*(1), 177–193.

Hotz, T., S. Huckemann, H. Le, J. S. Marron, J. Mattingly, E. Miller, J. Nolen, M. Owen, V. Patrangenaru, and S. Skwerer (2013). Sticky central limit theorems on open books. *Annals of Applied Probability 23*(6), 2238–2258.

Huckemann, S., J. C. Mattingly, E. Miller, and J. Nolen (2015). Sticky central limit theorems at isolated hyperbolic planar singularities. *Electronic Journal of Probability 20*(78), 1–34.

Huckemann, S. F. and B. Eltzner (2020). Data analysis on nonstandard spaces. *Wiley Interdisciplinary Reviews: Computational Statistics*, e1526.

Huckemann, S. F. and B. Eltzner (2024). *Foundations of Non-Euclidean Statistics.* London: Chapman & Hall/CRC Press. in preparation.

Hundrieser, S., B. Eltzner, and S. F. Huckemann (2020). Finite sample smeariness of Fréchet means and application to climate. *arXiv preprint arXiv:2005.02321*.

Pennec, X. (2019). Curvature effects on the empirical mean in Riemannian and affine manifolds: a non-asymptotic high concentration expansion in the small-sample regime. *arXiv preprint arXiv:1906.07418*.

Shevtsova, I. (2011). On the absolute constants in the Berry-Esseen type inequalities for identically distributed summands.

Sturm, K. (2003). Probability measures on metric spaces of nonpositive curvature. *Contemporary mathematics 338*, 357–390.

Tran, D., B. Eltzner, and S. F. Huckemann (2021). Smeariness begets finite sample smeariness. In *Geometric Science of Information 2021 proceedings*, pp. 29–36. Springer.