

Estimation of high dimensional gamma convolutions through random projections.

Thorin measures and Grassmannians cubatures...

O. Laverny ^{1,2}

April 20, 2021

¹ Institut Camille Jordan, UMR 5208, Université Claude Bernard Lyon 1, Lyon, France

² SCOR SE

Table of contents

1. Multivariate generalized gamma convolutions
2. Shifted moments, shifted cumulants and Laguerre coefficients
3. Direct minimisation of the Laguerre loss.
4. Better approach via random projections.
5. Examples

Multivariate generalized gamma convolutions

Construction of the multivariate Thorin class

Consider that X is a univariate gamma distribution with shape $\alpha \in \mathbb{R}_+$ and scale $s \in \mathbb{R}_+$. By definition, the moment generating function of X is

$$M(t) := \mathbb{E} \left(e^{tX} \right) = (1 - ts)^{-\alpha}.$$

Construction of the multivariate Thorin class

Consider that X is a univariate gamma distribution with shape $\alpha \in \mathbb{R}_+$ and scale $s \in \mathbb{R}_+$. By definition, the **cumulant** generating function of X is

$$K(t) := \ln \mathbb{E} \left(e^{tX} \right) = -\alpha \ln (1 - ts).$$

Construction of the multivariate Thorin class

Consider that \mathbf{X} is a **multivariate** gamma distribution with shape $\alpha \in \mathbb{R}_+$ and **scales** $\mathbf{s} \in \mathbb{R}_+^d$. By definition, the cumulant generating function of \mathbf{X} is

$$K(\mathbf{t}) := \ln \mathbb{E} \left(e^{\langle \mathbf{t}, \mathbf{X} \rangle} \right) = -\alpha \ln (1 - \langle \mathbf{t}, \mathbf{s} \rangle).$$

Warn: This distribution is comonotonous. We have $\mathbf{X} = (s_1 X, \dots, s_d X)$ with X a gamma distribution with shape α and unit scale.

Construction of the multivariate Thorin class

Consider that \mathbf{X} is a multivariate gamma convolution with shapes $\alpha \in \mathbb{R}_+^n$ and scales $\mathbf{s} \in \mathbb{R}_+^{n \times d}$. By definition, the cumulant generating function of \mathbf{X} is

$$K(\mathbf{t}) := \ln \mathbb{E} \left(e^{\langle \mathbf{t}, \mathbf{X} \rangle} \right) = - \sum_{i=1}^n \alpha_i \ln (1 - \langle \mathbf{t}, \mathbf{s}_i \rangle).$$

Warn: This distribution can be absolutely continuous w.r.t. λ .

Construction of the multivariate Thorin class

Consider that \mathbf{X} is a multivariate **generalized** gamma convolution with **Thorin measure** $\nu \in \mathcal{M}_+(\mathbb{R}_+^d)$. By definition, the cumulant generating function of X is

$$K(\mathbf{t}) := \ln \mathbb{E} \left(e^{\langle \mathbf{t}, \mathbf{X} \rangle} \right) = - \int \ln(1 - \langle \mathbf{t}, \mathbf{s}_i \rangle) \nu(\partial \mathbf{s}).$$

Good: This distribution can also be absolutely continuous w.r.t. λ , *under mild integration conditions*¹ on ν

We denote $\mathbf{X} \sim \mathcal{G}_d(\nu)$, and \mathcal{G}_d is called the d-variate Thorin class.

¹Victor Pérez-Abreu and Robert Stelzer. "A Class of Infinitely Divisible Multivariate and Matrix Gamma Distributions and Cone-Valued Generalised Gamma Convolutions". en. In: *arXiv:1201.1461 [math, stat]* (Jan. 2012).

The d -variate Thorin class \mathcal{G}_d

Definition (Multivariate Thorin Classes²)

$$\forall \nu \in \mathcal{M}_+(\mathbb{R}_+^d), \mathbf{X} \sim \mathcal{G}_d(\nu) \Leftrightarrow K(\mathbf{t}) := \ln \mathbb{E} (e^{\langle \mathbf{t}, \mathbf{X} \rangle}) = - \int \ln (1 - \langle \mathbf{s}, \mathbf{t} \rangle) \nu(d\mathbf{s}).$$

Prop: \mathcal{G}_d is closed w.r.t (independent) sums and products, and contains many interesting marginals...

²Lennart Bondesson. "On Univariate and Bivariate Generalized Gamma Convolutions". en. In: *Journal of Statistical Planning and Inference* 139.11 (Nov. 2009), pp. 3759–3765. ISSN: 03783758.

A convolutive representation (and motivations)

Note that in the finitely atomic case, say $\mathbf{X} \sim \mathcal{G}_d(\sum_{i=1}^n \alpha_i \delta_{\mathbf{s}_i})$, for some $\alpha \in \mathbb{R}_+^n$ and $\mathbf{s} \in \mathbb{R}_+^{n \times d}$, there exists independent Gamma random variables $G_i \sim \mathcal{G}_1(\alpha_i \delta_1)$, all having unit scale, such that:

$$\begin{pmatrix} X_1 \\ \dots \\ X_d \end{pmatrix} = \begin{pmatrix} s_{1,1} & \dots & \dots & s_{1,n} \\ \dots & \dots & \dots & \dots \\ s_{d,1} & \dots & \dots & s_{d,n} \end{pmatrix} \cdot \begin{pmatrix} G_1 \\ \dots \\ G_n \end{pmatrix},$$

Of course, we mostly consider $n \gg d$ (underdetermined case) with a sparse \mathbf{s} matrix.

Goal: Estimate ν from observations of the random vector \mathbf{X} .

Pb: It is a deconvolution problem, which is numerically hard.

Motivation: Risk factor interpretation and need for infinite divisibility...

Shifted moments, shifted cumulants and Laguerre coefficients

Estimation idea: orthonormal expansion of the density.

The idea is the following:

- (i) Find a **suitable** orthonormal basis of \mathbb{R}_+^d
- (ii) Expand the density into this basis, with an appropriate truncature.
- (iii) Compare theoretical and empirical coefficients to fit the parameters.

Bingo: The Laguerre basis³ provides usable closed form expression for many usefull quantities here.

³Florian Dussap. "Anisotropic multivariate deconvolution using projection on the Laguerre basis". In: *Journal of Statistical Planning and Inference* 215 (2021), pp. 23–46.

Laguerre coefficients \mathbf{a} and shifted moments $\boldsymbol{\mu}$

Definition (Laguerre basis of $L_2(\mathbb{R}_+^d)$)

$$\forall \mathbf{p} \in \mathbb{N}^d, \varphi_{\mathbf{p}}(\mathbf{x}) = \sqrt{2^d} \sum_{\mathbf{k} \leq \mathbf{p}} \binom{\mathbf{p}}{\mathbf{k}} \frac{(-2\mathbf{x})^{\mathbf{k}}}{\mathbf{k}!} e^{-|\mathbf{x}|}; \quad f(\mathbf{x}) = \sum_{\mathbf{k} \in \mathbb{N}^d} a_{\mathbf{k}} \varphi_{\mathbf{k}}(\mathbf{x})$$

$$\implies a_{\mathbf{k}} = \langle \varphi_{\mathbf{k}}, f \rangle = \mathbb{E}(\varphi_{\mathbf{k}}(\mathbf{X})) = \sqrt{2^d} \sum_{\mathbf{k} \leq \mathbf{p}} \binom{\mathbf{p}}{\mathbf{k}} \frac{(-2)^{|\mathbf{k}|}}{\mathbf{k}!} \mathbb{E}\left(\mathbf{x}^{\mathbf{k}} e^{-|\mathbf{x}|}\right)$$

Denote $\mu_i := \mathbb{E}\left(\mathbf{x}^{\mathbf{k}} e^{-|\mathbf{x}|}\right)$. Let $I \subset \mathbb{N}^d$ an increasing index set, $\mathbf{a} = (a_i)_{i \in I}$ and $\boldsymbol{\mu} = (\mu_i)_{i \in I}$.

Bijection: The relationship between \mathbf{a} and $\boldsymbol{\mu}$ is linear. We encode this relationship through a (lower-triangular, invertible) matrix $\mathbf{A} \in \mathbb{R}^{|I| \times |I|}$ such that

$$\mathbf{a} = \mathbf{A}\boldsymbol{\mu} \text{ and } \boldsymbol{\mu} = \mathbf{A}^{-1}\mathbf{a}.$$

Thorin moments τ

Rem: Since $\mu_i := \mathbb{E}(\mathbf{X}^{\mathbf{k}} e^{-\langle \mathbf{1}, \mathbf{X} \rangle})$, the mgf of the random vector writes:

$$M(\mathbf{t}) = \mathbb{E}(e^{\langle \mathbf{t}, \mathbf{X} \rangle}) = \sum_{\mathbf{k} \in \mathbb{N}^d} \mu_{\mathbf{k}} \frac{(\mathbf{t} - \mathbf{1})^{\mathbf{k}}}{\mathbf{k}!},$$

Definition (Thorin moments τ)

$$K(\mathbf{t}) = \ln M(\mathbf{t}) := \sum_{\mathbf{k} \in \mathbb{N}^d} \tau_{\mathbf{k}} (|\mathbf{k}| - 1)! \frac{(\mathbf{t} - \mathbf{1})^{\mathbf{k}}}{\mathbf{k}!}.$$

Bijection: There exists a function \mathbf{B} , based on Bell polynomials, s.t.

$$\mu = \mathbf{B}(\tau) \text{ and } \tau = \mathbf{B}^{-1}(\mu).$$

Estimators of a, μ and τ

Data: $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathbb{R}_+^{N \times d}$ an N -sample of i.i.d. random vectors.

Definition (Monte-Carlo estimators)

$$\hat{\mu}(\mathbf{x}) = (\widehat{\mu}_k(\mathbf{x}))_{k \in I} = \left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^k e^{-|\mathbf{x}_i|} \right)_{k \in I}$$

$$\hat{\tau}(\mathbf{x}) = (\widehat{\tau}_k(\mathbf{x}))_{k \in I} = \mathbf{B}^{-1}(\hat{\mu}(\mathbf{x})).$$

$$\hat{a}(\mathbf{x}) = (\widehat{a}_k(\mathbf{x}))_{k \in I} = \mathbf{A}\hat{\mu}(\mathbf{x}) = \mathbf{A}\mathbf{B}(\hat{\tau}(\mathbf{x})),$$

Rem: biais, cv ?

Back to gamma convolutions...

Definition (Thorin moments of a Gamma convolution)

We denote the first Thorin moments of the $\mathcal{G}_d(\nu)$ distribution by:

$$\tau(\nu) = (\tau_k(\nu))_{k \in I}.$$

Property

Denoting δ_x the Dirac measure at x , it holds:

- (i) $\tau_0(\delta_s) = -\ln(1 + |s|)$ and $\tau_k(\delta_s) = \left(\frac{s}{1+|s|}\right)^{|k|}$ for $k \neq 0$, $k \in \mathbb{N}^d$.
- (ii) $\tau(\nu) = \int \tau(\delta_s) \nu(ds)$, where the equality and integration are intended component-wise.

Rem: τ are called Thorin moments for a reason. We are dealing with a multivariate moment problem *which might not have a solution...*

Direct minimisation of the Laguerre loss.

We use the integrated square error between densities, projected in the Laguerre basis:

$$\mathcal{L}(\mathbf{x}, \nu) = \|\hat{\mathbf{a}}(\mathbf{x}) - \mathbf{AB}(\tau(\nu))\|_2^2.$$

Theorem (Consistency⁴)

If \mathbf{x} is drawn from an ϵ -well-behaved density $f \in \mathcal{G}_d$, any well-behaved estimator ν^* such that $\mathcal{L}(\mathbf{x}, \nu^*) \xrightarrow[N \rightarrow \infty]{a.s.} 0$ ensures that

$$\|f - f_{\nu^*}\|_2^2 \xrightarrow[N \rightarrow \infty, I \rightarrow \mathbb{N}^d]{a.s.} 0.$$

⁴Oskar Laverny, Esterina Masiello, Véronique Maume-Deschamps, and Didier Rulli  re. "Estimation of multivariate generalized gamma convolutions through Laguerre expansions.". In: *Electronic Journal of Statistics* 15.2 (2021), pp. 5158–5202. DOI: 10.1214/21-EJS1918. URL: <https://doi.org/10.1214/21-EJS1918>.

The loss: $\mathcal{L}(\mathbf{x}, \nu) = \|\hat{\mathbf{a}}(\mathbf{x}) - \mathbf{AB}(\tau(\nu))\|_2^2$.

Problem: \mathcal{L} is too costly to work with when d gets large...

The vectors $\hat{\mathbf{a}}(\mathbf{x})$ and $\mathbf{AB}(\tau(\nu))$ each consist of $|I|$ coefficients. If, e.g., $I = \{\mathbf{k} \in \mathbb{N}^d, |\mathbf{k}| \leq m\}$ is isotropic, the number of coefficients to compute is given by

$$D(m, d) = \sum_{i=1}^m \binom{i + d - 1}{d - 1}.$$

which is exponentially increasing in d and therefore unusable.

**Better approach via random
projections.**

An approximated loss through random projections...

Through a first-order taylor expansion of the function \mathbf{AB} , we define:

$$\widehat{\mathcal{L}}(\mathbf{x}, \nu) = \|\hat{\boldsymbol{\tau}}(\mathbf{x}) - \boldsymbol{\tau}(\nu)\|_{\nabla(\mathbf{x})}^2$$

where $\nabla(\mathbf{x})$ is a jacobian of the function \mathbf{AB} , taken in $\hat{\boldsymbol{\tau}}(\mathbf{x})$. Then, through a univariate projection and re-integration we define:

$$\begin{aligned}\widetilde{\mathcal{L}}(\mathbf{x}, \nu) &:= \int_{[0,1]^d} \widehat{\mathcal{L}}(\langle \mathbf{c}, \mathbf{x} \rangle, \nu_{\langle \mathbf{c} \rangle}) d\mathbf{c} \\ &= \int_{[0,1]^d} \|\hat{\boldsymbol{\tau}}(\langle \mathbf{c}, \mathbf{x} \rangle) - \boldsymbol{\tau}(\nu_{\langle \mathbf{c} \rangle})\|_{\nabla(\langle \mathbf{c}, \mathbf{x} \rangle)}^2 d\mathbf{c},\end{aligned}$$

where

$$\nu_{\langle \mathbf{c} \rangle}(A) = \nu\left(\left\{\mathbf{x} \in \mathbb{R}_+^d : \langle \mathbf{c}, \mathbf{x} \rangle \in A\right\}\right) \quad \forall A \subseteq \mathbb{R}_+.$$

... that is still consistent.

Theorem (Consistency of $\tilde{\mathcal{L}}$.)

Let \mathbf{x} be drawn from a well-behaved density $f \sim \mathcal{G}_d(\nu)$. The global minimizer

$$\tilde{\nu} := \arg \min_{\nu: \mathcal{G}_d(\nu) \text{ w.b.}} \tilde{\mathcal{L}}(\mathbf{x}, \nu)$$

ensures that

$$\|f - f_{\tilde{\nu}}\|_2^2 \xrightarrow[N \rightarrow \infty]{a.s.} 0.$$

The proof leverages the theory of Grassmannian cubatures and some considerations about unisolvant sets for polynomials of bounded degree.

Question: Can we and will we reach a global minimizer ? Recall that the loss is clearly not convex and has a myriad of local minimas...

Reframing the loss

Let $\mathcal{F} = L_2([0, 1]^d, \mathbb{R}_+^m)$ be the space of square integrable functions from $[0, 1]^d$ to \mathbb{R}_+^m .

We consider the functional ψ defined as:

$$\psi: \begin{cases} \mathbb{R}_+^d \longrightarrow \mathcal{F} \\ \mathbf{s} \longmapsto \psi(\mathbf{s}): \begin{cases} [0, 1]^d \longrightarrow \mathbb{R}_+^m \\ \mathbf{c} \longmapsto \tau(\delta_{\langle \mathbf{s}, \mathbf{c} \rangle}) . \end{cases} \end{cases}$$

We also consider, as our target,

$$\psi_{\mathbf{x}}: \begin{cases} [0, 1]^d \longrightarrow \mathbb{R}_+^m \\ \mathbf{c} \longmapsto \hat{\tau}(\langle \mathbf{c}, \mathbf{x} \rangle) \end{cases} \in \mathcal{F}.$$

Reframing the loss

We endow \mathcal{F} with the norm $\|\cdot\|_{\mathbf{x}}$ defined as:

$$\forall \psi \in \mathcal{F}, \|\psi\|_{\mathbf{x}}^2 = \int_{[0,1]^d} \|\psi(\mathbf{c})\|_{\nabla(\langle \mathbf{c}, \mathbf{x} \rangle)}^2 d\mathbf{c},$$

This allow us to write our loss as:

$$\tilde{\mathcal{L}}(\mathbf{x}, \nu) = \|\psi_{\mathbf{x}} - \int \psi(\mathbf{s}) \nu(d\mathbf{s})\|_{\mathbf{x}}^2,$$

which shows that our loss is convex in $\int \psi(\mathbf{s}) \nu(d\mathbf{s})$, but clearly not convex in the atoms and weights of an n -atomic measure ν . However Chizat⁵ applies !

⁵Lenaic Chizat. "Sparse optimization on measures with over-parameterized gradient descent". In: *Mathematical Programming* (2021), pp. 1–46.

Fréchet-differentiability and gradient flow.

Property

The Fréchet differential of $\tilde{\mathcal{L}}(\mathbf{x}, \cdot)$ at $\nu \in \mathcal{M}(\mathbb{R}_+^d)$ is represented by the function

$$\tilde{\mathcal{L}}'_\nu: \begin{cases} \mathbb{R}_+^d \longrightarrow \mathbb{R} \\ \mathbf{s} \longmapsto 2\langle \psi_{\mathbf{x}} - \int \psi(\mathbf{s})\nu(d\mathbf{s}), \psi(\mathbf{s}) \rangle_{\mathbf{x}}. \end{cases}$$

Definition (Gradient flow⁶)

A gradient flow of $\tilde{\mathcal{L}}(\mathbf{x}, \cdot)$ is an absolutely continuous curve $(\nu_t)_{t \geq 0}$ in the space $\mathcal{M}(\mathbb{R}_+^d)$ that satisfies

$$\frac{\partial}{\partial t} \nu_t = -\nabla \tilde{\mathcal{L}}(\mathbf{x}, \nu_t)$$

Where the gradient $\nabla \tilde{\mathcal{L}}(\mathbf{x}, \nu_t)$ is taken w.r.t. a given conic metric...

⁶Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2005.

A convergence result for the gradient flow.

Theorem (Global convergence of the gradient flow)

For $\rho \in \mathcal{M}_+(\mathbb{R}_+^d)$ an absolutely continuous reference measure such that $\log \rho$'s density is Lipschitz, for any initial measure $\nu_0 \in \mathcal{M}_+(\mathbb{R}_+^d)$, there exists a constant C , dependent on the characteristics of the problem, such that if $W_\infty(\nu_0, \rho) \leq C$,

$$\exists \nu_\infty \in \arg \min_{\nu \in \mathcal{M}_+(\mathbb{R}_+^d)} \tilde{\mathcal{L}}(\mathbf{x}, \nu) \text{ such that } W_\infty(\nu_t, \nu_\infty) \xrightarrow[t \rightarrow \infty]{} 0.$$

Furthermore, when $\nu_0 = \rho$, we achieve a precision ϵ , i.e., $W_\infty(\nu_t, \nu_\infty) \leq \epsilon$, provided the number of iteration is $t = \mathcal{O}(-\log(\epsilon))$.

This result directly leverages Chizat's results. For general accelerated convex methods, $W_\infty(\nu_t, \nu_\infty) \leq \epsilon$ is achieved for $t = \mathcal{O}(\epsilon^{-1/d})$ only, see⁷ for more details.

⁷Yohann de Castro, Sébastien Gadat, Clément Marteau, and Cathy Maugis. "SuperMix: Sparse Regularization for Mixtures". In: *The annals of Statistics* 49.3 (2021), pp. 1779–1809.

Takeaway: Initialize the Thorin measure all over \mathbb{R}_+^{d+1} , with a lot of atoms, to achieve global convergence.

Bonus: Same result with a lasso penalty on the measure, which is perfect for what we need (penalises the abundance of sources), via a hyperparameter $\lambda > 0$ and a loss

$$\tilde{\mathcal{L}}(\mathbf{x}, \nu) + \lambda |\nu|.$$

Any wanted hyperparameter searching method (with its associated cost) could be used to find λ .

Examples

Estimation procedure

Algorithm 1: Estimation of Thorin measures via stochastic gradient descent on $\tilde{\mathcal{L}}$.

Input: A dataset $\mathbf{x} \in \mathbb{R}^{N \times d}$, a number of Gammas $n \in \mathbb{N}$, a precision parameter $m \in \mathbb{N}$, a number of iterations $T \in \mathbb{N}$, and a learning rate $\eta \in \mathbb{R}_+$

Result: A Thorin measure ν_T that approximates the dataset \mathbf{x} as a multivariate Gamma convolution.

Estimate standard deviations $\sigma_i = \text{std}(\mathbf{x}_i)$ for all $i \in 1, \dots, d$, and standardize the marginals by dividing \mathbf{x}_i by σ_i .

Initialize a measure $\nu_0 \in \mathcal{M}_+(\mathbb{R}_+^d)$ with n atoms and corresponding weights, chosen randomly through Gaussian noise.

foreach $t \in 0, \dots, T - 1$ **do**

 Choose a random direction $\mathbf{c} \in [0, 1]^d$.

 Compute the gradient \mathbf{g} of $\hat{\mathcal{L}}(\langle \mathbf{c}, \mathbf{x} \rangle, \nu_{t, \langle \mathbf{c} \rangle})$ with respect to ν_t (details missing).

 Let $\nu_{t+1} = \nu_t - \eta \mathbf{g}$

end

Rescale ν_T by $\sigma_1, \dots, \sigma_d$.

Return ν_T

Four-variate simulated data

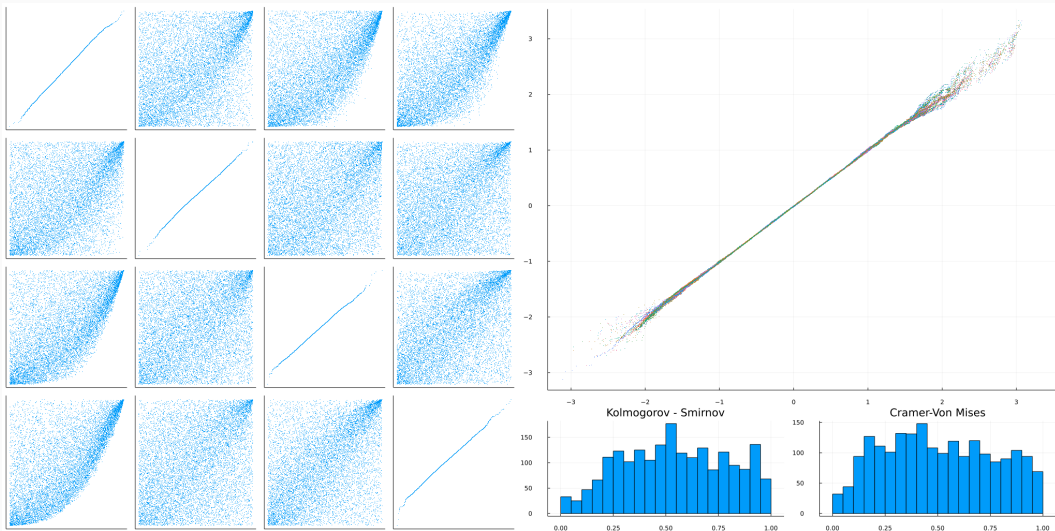
For U_1, U_2, U_3 independent $\mathcal{U}([0, 1])$ distributions and Y_1, \dots, Y_4 independent $\text{log-Normal}(0, 1)$ distributions, we let the random vector \mathbf{X} be defined by:

$$\mathbf{X} = \left(Y_1, Y_2 + U_1 Y_1^2, Y_3 + U_3 Y_1, Y_4 + Y_1^{1 + \frac{U_3}{3}} \right),$$

We simulate a dataset $\mathbf{x} \in \mathbb{R}^{10000 \times 4}$ of $N = 10000$ i.i.d. samples from \mathbf{X} , and we ran the algorithm on it with 100 atoms.

We ended up with only 17 atoms at a $1e - 16$ threshold on weights.

Four-variate simulated data



2000-dimensional multiplicative dataset

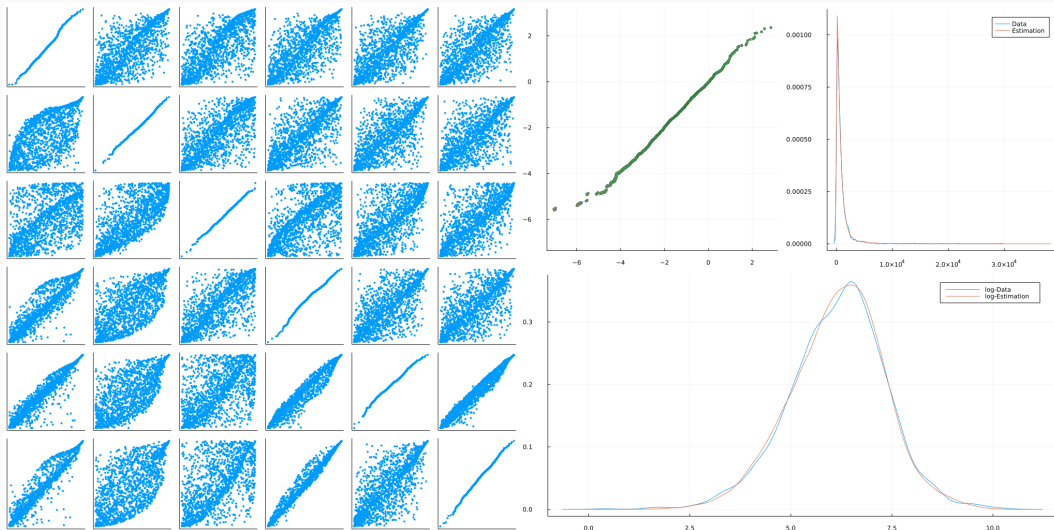
Let $d = 2000$, let $G \sim \mathcal{G}_1(\delta_1)$, $H \sim \mathcal{G}_1(2\delta_{\frac{1}{2}})$, and $Z_1, \dots, Z_d \sim \mathcal{N}(0, 1)$ be all independent random variables, and let $\alpha_1, \dots, \alpha_d$ be fixed parameters in $[0, 1]$ (uniformly chosen). Construct the random vector $\mathbf{X} = (X_1, \dots, X_d)$ as:

$$\mathbf{X} = \left(G e^{Z_i} H^{1+2\alpha_i} \right)_{i \in 1, \dots, d}.$$

We simulate a dataset $\mathbf{x} \in \mathbb{R}^{1500 \times 2000}$ of $N = 1500$ i.i.d. samples from \mathbf{X} .

We ran our algorithm with $n = 500$ atoms, and ended up with only 36 meaningful ones.

2000-dimensional multiplicative dataset



Conclusion

Main takeaways:

- (i) Multivariate generalized gamma convolutions are flexible semi-parametric structures
- (ii) They simplify divisions of positive random variables (in the infinite divisibility sense) by making this process parametrical
- (iii) Estimating them can be reduced to a non-sparse d -variate moment problem, which can be very hard to solve.
- (iv) Random projections allow us to make the gradient cost essentially linear in the dimension, and still converges to a globally minimizing Thorin measure.
- (v) We achieve sparse Thorin measure from dense proposals.
- (vi) The (OSS) Julia package `ThorinDistributions.jl` provides an implementation.

Thanks !