

Some intriguing properties of extreme geometric quantiles

Stéphane GIRARD (Inria Grenoble Rhône-Alpes)
Joint work with Gilles STUPFLER (Université Aix-Marseille)

Paris, January 2015

Outline

- Geometric quantiles
- Extreme geometric quantiles
 - Under moment conditions
 - In a multivariate regular variation framework
- Numerical illustrations
- Real data example

Multivariate quantiles

- The natural order on \mathbb{R} induces a universal definition of quantiles for univariate distribution functions.
- This is not true in \mathbb{R}^d , $d \geq 2$, since no natural order exists in this case.
- Many definitions of multivariate quantiles have thus been suggested in the literature:
 - **Depth-based quantiles:** Liu *et al.* (1999), Zuo & Serfling (2000)
 - **Norm minimisation:** Abdous & Theodorescu (1992), Chaudhuri (1996)

For a review, see e.g. Serfling (2002).

Geometric quantiles

If X is a real-valued random variable, its univariate p -th quantile

$$q(p) := \inf\{t \in \mathbb{R} \text{ s.t. } \mathbb{P}(X \leq t) \geq p\}$$

can be obtained by solving the optimisation problem

$$\arg \min_{q \in \mathbb{R}} \mathbb{E}(|X - q| - |X|) - (2p - 1)q.$$

- When $\mathbb{E}|X| < \infty$, this problem can be simplified as

$$\arg \min_{q \in \mathbb{R}} \mathbb{E}|X - q| - (2p - 1)q.$$

In particular, the **median** $q(1/2)$ of X is obtained by minimising $\mathbb{E}|X - q|$ with respect to q .

- Subtracting $\mathbb{E}|X|$ makes the cost function well-defined even when $\mathbb{E}|X| = \infty$.

In \mathbb{R}^d , $d \geq 2$, analogues of the absolute value and product are given by the Euclidean norm $\|\cdot\|$ and Euclidean scalar product $\langle \cdot, \cdot \rangle$.

When X is a multivariate random vector, the **geometric quantiles** of X , introduced by Chaudhuri (1996), are thus obtained by adapting the aforementioned problem in the multivariate context:

Definition 1 (Chaudhuri 1996)

If $u \in \mathbb{R}^d$ is an arbitrary vector, a geometric u -th quantile of X , if it exists, is a **solution of the optimisation problem**

$$\arg \min_{q \in \mathbb{R}^d} \mathbb{E}(\|X - q\| - \|X\|) - \langle u, q \rangle. \quad (P_u)$$

Properties

Central properties

- For all $u \in \mathbb{R}^d$ such that $\|u\| < 1$, there exists a unique geometric u -th quantile whenever the distribution of X is not concentrated on a single straight line in \mathbb{R}^d (Chaudhuri, 1996).
- They are equivariant under any orthogonal transformation (Chaudhuri, 1996).
- The geometric quantile function characterises the associated distribution (Koltchinskii, 1997).

These central properties make geometric quantiles reasonable candidates when trying to define multivariate quantiles.

Extreme properties? Our focus here is to investigate the properties of extreme geometric quantiles.

A first step

From now on, we assume that the distribution of X is not concentrated on a single straight line in \mathbb{R}^d and non-atomic. Then:

Proposition 1

The optimisation problem (P_u) has a solution if and only if $\|u\| < 1$.

- We cannot compute a geometric quantile with unit index vector, unlike in the univariate case when the distribution has a finite (left or right) endpoint.
- We may nevertheless study the asymptotics of a geometric quantile $q(u)$ when $\|u\| \rightarrow 1$: such quantiles will be referred to as **extreme geometric quantiles** (Chaudhuri, 1996, Cheng and De Gooijer, 2007) and (Chaouch and Goga, 2010) for outlier detection.

Theorem 1

Let $u \in \mathbb{R}^d$.

- (i) $\|q(u)\| \rightarrow \infty$ as $\|u\| \rightarrow 1$.
- (ii) Moreover, if $\|u\| = 1$ and $\lambda \uparrow 1$ then

$$q(\lambda u) / \|q(\lambda u)\| \rightarrow u.$$

- The magnitude of extreme geometric quantiles diverges to infinity. (Rather intriguing: it holds true even if the distribution of X has a compact support. A related point: sample geometric quantiles do not necessarily lie within the convex hull of the sample (Breckling *et al.* 2001)).
- If $\|u\| = 1$ and $\lambda \uparrow 1$ then the extreme geometric quantile $q(\lambda u)$ has asymptotic direction u .

The next results specify rates of the convergence in Theorem 1 under further assumptions.

When there are finite moments

Our first result is obtained in the case when $\|X\|$ satisfies certain moment conditions. It focuses on extreme geometric quantiles in the direction u , i.e. having the form $q(\lambda u)$, with $\lambda \uparrow 1$ and $\|u\| = 1$.

Theorem 2

Let $u \in \mathbb{R}^d$ such that $\|u\| = 1$. Define $\Pi_u(x) = x - \langle x, u \rangle u$.

(i) If $\mathbb{E}\|X\| < \infty$ then

$$\|q(\lambda u)\| \left(\frac{q(\lambda u)}{\|q(\lambda u)\|} - u \right) \rightarrow \mathbb{E}(\Pi_u(X)) \text{ as } \lambda \uparrow 1.$$

(ii) If $\mathbb{E}\|X\|^2 < \infty$ and Σ denotes the *covariance matrix* of X then

$$\|q(\lambda u)\|^2(1 - \lambda) \rightarrow \frac{1}{2} (\text{tr } \Sigma - u' \Sigma u) > 0 \text{ as } \lambda \uparrow 1.$$

Consequences of Theorem 2

If $\|X\|$ has a finite second moment, then:

- The asymptotic direction of an extreme geometric quantile in the direction u is exactly u .
- The magnitude of an extreme geometric quantile in the direction u is asymptotically determined by u and the covariance matrix Σ .

In particular, the extreme geometric quantiles of two probability distributions with the same finite covariance matrices are asymptotically equivalent.

⇒ No information can be recovered on the “tail” behaviour of the distribution basing solely on extreme geometric quantiles.

Further consequences of Theorem 2

If $\|X\|$ has a finite second moment, then:

- Shape of an extreme quantile contour: The global maximum of the function $u \mapsto \text{tr} \Sigma - u' \Sigma u$ on the unit sphere is reached at a unit eigenvector of Σ associated with its smallest eigenvalue. Thus, the norm of an extreme geometric quantile is the largest in the direction where the variance is the smallest.
- It is possible to estimate an extreme geometric quantile parametrically:

$$\hat{q}_n(\alpha_n u) = (1 - \alpha_n)^{-1/2} \left[\frac{1}{2} \left(\text{tr} \hat{\Sigma}_n - u' \hat{\Sigma}_n u \right) \right]^{1/2} u.$$

where $\hat{\Sigma}_n$ is the empirical counterpart of Σ . It is very different from the univariate case which requires semi-parametric estimators of tail indices.

In a multivariate regular variation framework

When the moment conditions are no longer satisfied, the asymptotic properties of extreme geometric quantiles can be studied in a **multivariate regular variation framework**:

(M_α) The random vector X has a probability density function f which is continuous on a neighborhood of infinity and such that:

- There exist a positive function Q on \mathbb{R}^d and a function V which is **regularly varying at infinity** with index $-\alpha < 0$, such that

$$\forall y \neq 0, \left| \frac{f(ty)}{t^{-d}V(t)} - Q(y) \right| \rightarrow 0 \text{ as } t \rightarrow \infty$$

and $\sup_{\|w\|=1} \left| \frac{f(tw)}{t^{-d}V(t)} - Q(w) \right| \rightarrow 0 \text{ as } t \rightarrow \infty.$

- The function $y \mapsto \|y\|^d f(y)$ is locally bounded at 0.

This model is closely related to the one of Cai *et al.* (2011).

If (M_α) holds, then:

- The function Q is **homogeneous** of degree $-d - \alpha$ on $\mathbb{R}^d \setminus \{0\}$.
- We have that

$$f(x) = \|x\|^{-d} V(\|x\|) Q(x/\|x\|) (1 + o(1)),$$

as $\|x\| \rightarrow \infty$ and thus $f(x)$ is roughly of order $\|x\|^{-d-\alpha}$.

- The expectation $\mathbb{E}\|X\|^\beta$ is finite if $\beta < \alpha$.

In particular, the cases $\alpha > 2$ and $\alpha > 1$ are covered by Theorem 2.

Theorem 3

Let $u \in \mathbb{R}^d$ such that $\|u\| = 1$.

(i) If (M_α) holds with $\alpha \in (0, 1)$, then

$$\frac{1}{V(\|q(\lambda u)\|)} \left(\frac{q(\lambda u)}{\|q(\lambda u)\|} - u \right) \rightarrow \int_{\mathbb{R}^d} \frac{\Pi_u(y)}{\|y - u\|} Q(y) dy \quad \text{as } \lambda \uparrow 1.$$

(ii) If (M_α) holds with $\alpha \in (0, 2)$, then

$$\frac{1 - \lambda}{V(\|q(\lambda u)\|)} \rightarrow \int_{\mathbb{R}^d} \left(1 + \frac{\langle y - u, u \rangle}{\|y - u\|} \right) Q(y) dy \quad \text{as } \lambda \uparrow 1.$$

Comments on Theorem 3

Since V is regularly varying with index $-\alpha$, it follows that when $\alpha \in (0, 2)$, the magnitude of an extreme geometric quantile roughly behaves like $(1 - \lambda)^{-1/\alpha}$ as $\lambda \uparrow 1$.

\Rightarrow In this case, the magnitude of an extreme geometric quantile features the “tail” behaviour of the distribution.

However, Theorem 3 excludes the limit cases $\alpha = 1$ for the asymptotic direction and $\alpha = 2$ for the asymptotic magnitude. These limit cases can be studied via sub-models of (M_α) .

To summarize ...

For all $\alpha > 0$, we can write

$$\frac{q(\lambda u)}{\|q(\lambda u)\|} - u \propto R_{1,\alpha}((1-\lambda)^{-1})$$

and $\|q(\lambda u)\| \propto R_{2,\alpha}((1-\lambda)^{-1})$ as $\lambda \uparrow 1$,

where $R_{1,\alpha}$ and $R_{2,\alpha}$ are **regularly varying functions** with respective indices $-\min(1, \alpha)/\min(2, \alpha)$ and $1/\min(2, \alpha)$.

\Rightarrow Extreme geometric quantiles feature the “tail” behaviour of X only when the density of $\|X\|$ decays sufficiently slowly at infinity.

Numerical illustrations: Theorem 2

We choose $d = 2$ to make the display easier. The following two bivariate distributions are considered:

- the centred Gaussian bivariate distribution $\mathcal{N}(0, v_X, v_Y, v_{XY})$ with covariance matrix

$$\Sigma = \begin{pmatrix} v_X & v_{XY} \\ v_{XY} & v_Y \end{pmatrix}.$$

- a centred double exponential distribution $\mathcal{E}(\lambda_-, \mu_-, \lambda_+, \mu_+)$, with $\lambda_-, \mu_-, \lambda_+, \mu_+ > 0$, whose probability density function is:

$$f(x, y) = \begin{cases} \frac{\lambda_+ \mu_+}{4} e^{-\lambda_+ |x| - \mu_+ |y|} & \text{if } xy > 0, \\ \frac{\lambda_- \mu_-}{4} e^{-\lambda_- |x| - \mu_- |y|} & \text{if } xy \leq 0. \end{cases}$$

In this case, X has an explicit covariance matrix $\Sigma(\lambda_-, \mu_-, \lambda_+, \mu_+)$.

Since both distributions have a finite covariance matrix, Theorem 2 entails that their extreme geometric quantiles are **asymptotically** equal to:

$$q_{\text{eq}}(\lambda u) := (1 - \lambda)^{-1/2} \left[\frac{1}{2} (\text{tr} \Sigma - u' \Sigma u) \right]^{1/2} u.$$

⇒ **Goal**: to show that for these two distributions, equal covariance matrices induce equivalent extreme geometric quantiles, and to assess the accuracy of the asymptotic equivalent.

We choose three different sets of parameters, in order that the related covariance matrices coincide:

- $\mathcal{N}(0, 1/2, 1/2, 0)$ and $\mathcal{E}(2, 2, 2, 2)$ with **spherical** covariance matrices;
- $\mathcal{N}(0, 1/8, 3/4, 0)$ and $\mathcal{E}(4, 2\sqrt{2/3}, 4, 2\sqrt{2/3})$ with **diagonal** but **non-spherical** covariance matrices;
- $\mathcal{N}(0, 1/2, 1/2, 1/6)$ and $\mathcal{E}(2\sqrt{3}, 2\sqrt{3}, 2\sqrt{3/5}, 2\sqrt{3/5})$ with **full** covariance matrices.

Any unit vector u can be written $u = u_\theta = (\cos \theta, \sin \theta)$, $\theta \in [0, 2\pi)$.

We let $\lambda = 0.995$ and in each case, we compute:

- the **true iso-quantile curve** $\mathcal{C}q(\lambda) = \{q(\lambda u_\theta), \theta \in [0, 2\pi)\}$;
- its **asymptotic equivalent** $\mathcal{C}q_{\text{eq}}(\lambda) = \{q_{\text{eq}}(\lambda u_\theta), \theta \in [0, 2\pi)\}$.

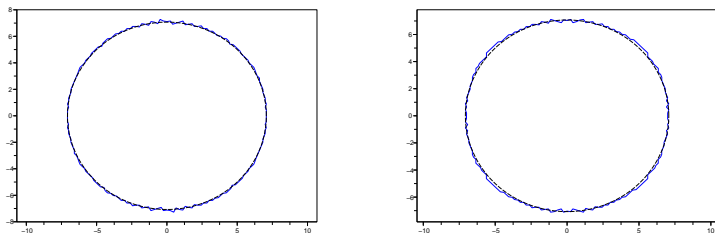


Figure 1: Spherical case. Gaussian (left) and double exponential (right) distributions. Iso-quantile curves $Cq(\lambda)$ (full blue line) and $Cq_{eq}(\lambda)$ (dashed black line).

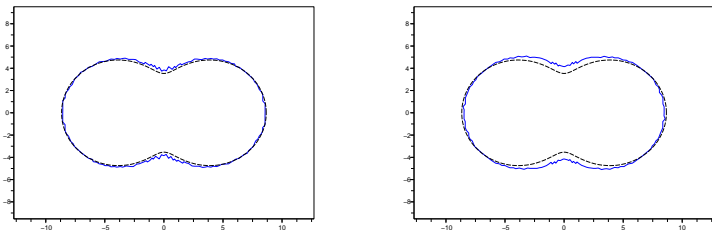


Figure 2: Diagonal case. Gaussian (left) and double exponential (right) distributions. Iso-quantile curves $C_q(\lambda)$ (full blue line) and $C_{q_{eq}}(\lambda)$ (dashed black line).

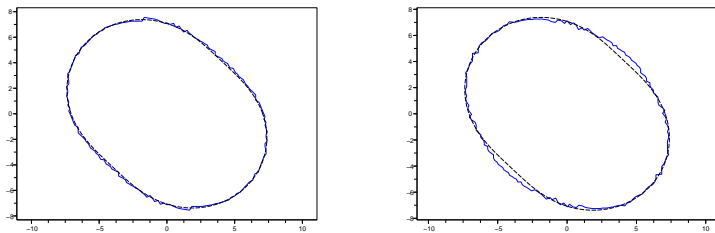


Figure 3: Full case. Gaussian (left) and double exponential (right) distributions. Iso-quantile curves $\mathcal{C}q(\lambda)$ (full blue line) and $\mathcal{C}q_{eq}(\lambda)$ (dashed black line).

Numerical illustrations: Theorem 3

Here, we consider a bivariate Pareto($\alpha, \sigma_1, \sigma_2$) distribution, whose probability density function is:

$$f(x, y) = \frac{\alpha}{2\sigma_1\sigma_2\pi} \left(\frac{x^2}{\sigma_1^2} + \frac{y^2}{\sigma_2^2} \right)^{(-2-\alpha)/2} \mathbb{1}_{[1, \infty)} \left(\frac{x^2}{\sigma_1^2} + \frac{y^2}{\sigma_2^2} \right)$$

where α , σ_1^2 and $\sigma_2^2 > 0$. When $\alpha > 2$, this distribution has covariance matrix:

$$\frac{1}{2} \cdot \frac{\alpha}{\alpha - 2} \Sigma, \quad \text{with } \Sigma = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}.$$

Clearly, for any $\alpha > 0$, this distribution is part of the class (M_α) , with

$$Q(x) = (x' \Sigma^{-1} x)^{(-2-\alpha)/2}$$

and $V(t) = \frac{\alpha}{2\sigma_1\sigma_2\pi} t^{-\alpha} \mathbb{1}_{[1,\infty)}(t).$

Theorems 2 and 3 thus entail that the extreme geometric quantiles of this distribution are **asymptotically** equal to:

$$q_{\text{eq}}(\lambda u) := (1 - \lambda)^{-1/\alpha} l(\alpha, \sigma_1, \sigma_2) u \quad \text{if } \alpha < 2$$

where $l(\alpha, \sigma_1, \sigma_2)$ is a positive constant, and

$$q_{\text{eq}}(\lambda u) := (1 - \lambda)^{-1/2} \left[\frac{1}{2} (\text{tr } M - u' M u) \right]^{1/2} u \quad \text{if } \alpha > 2.$$

⇒ **Goal**: to examine if both these approximations are satisfactory on this heavy-tailed example.

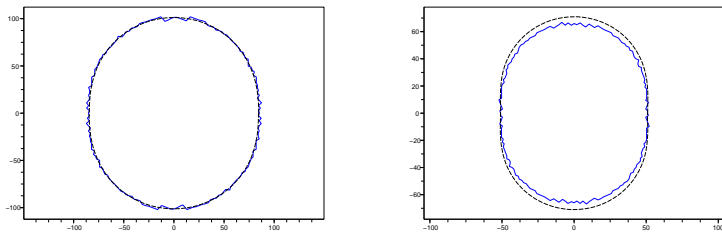


Figure 4: Pareto($\alpha, 2, 1/2$) model, with $\alpha = 1.3$ (left) and $\alpha = 1.5$ (right). Iso-quantile curves $\mathcal{C}q(\lambda)$ (full blue line) and $\mathcal{C}q_{\text{eq}}(\lambda)$ (black dashed line).

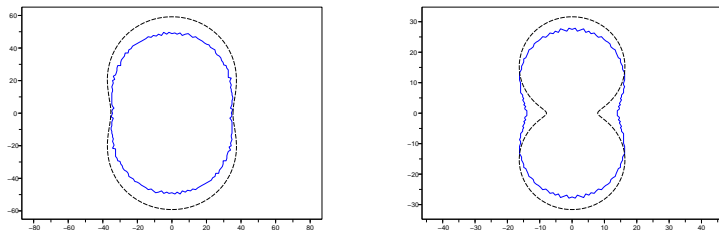


Figure 5: Pareto($\alpha, 2, 1/2$) model, with $\alpha = 1.7$ (left) and $\alpha = 2.5$ (right). Iso-quantile curves $\mathcal{C}q(\lambda)$ (full blue line) and $\mathcal{C}q_{eq}(\lambda)$ (black dashed line).

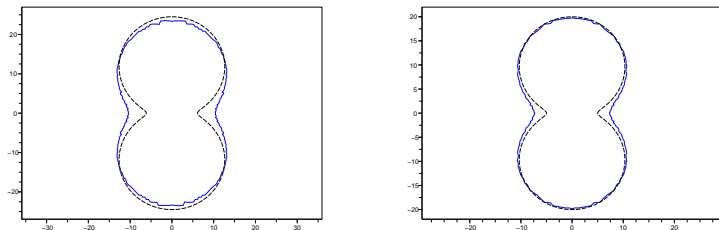


Figure 6: Pareto($\alpha, 2, 1/2$) model, with $\alpha = 3$ (left) and $\alpha = 4$ (right). Iso-quantile curves $Cq(\lambda)$ (full blue line) and $Cq_{eq}(\lambda)$ (black dashed line).

Illustration on the Pima Indians Diabetes Database

- The sample behaviour of extreme geometric quantiles is illustrated on a two-dimensional dataset extracted from the Pima Indians Diabetes Database

`ftp.ics.uci.edu/pub/machine-learning-databases/pima-indians-diabetes.`

- The data set consists of $n = 392$ pairs (X_i, Y_i) , where X_i is the body mass index of the i th individual and Y_i is its diastolic blood pressure.

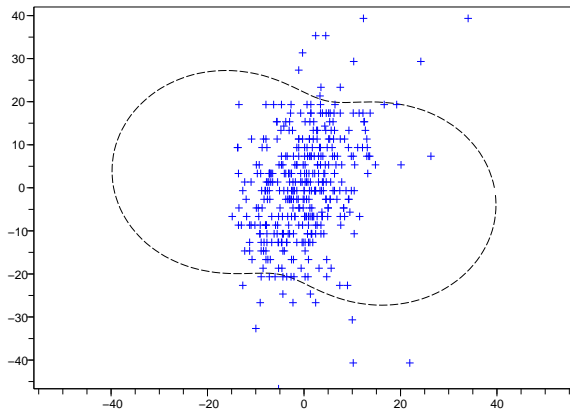


Figure 7: Centered data and estimated geometric iso-quantile curve, $\alpha = 0.95$.

Discussion

- Extreme geometric quantiles in the direction u have asymptotic direction u .
- They are asymptotically equal for two distributions which have the same finite covariance matrix, which is not satisfying from the extreme value perspective.
- The shape of the iso-quantile curves may be totally different from the shape of the density contour plots. Outlier detection should be conducted with great care.
- They do however feature the “tail” behaviour of X in a multivariate regular variation context when the tail of $\|X\|$ is sufficiently heavy.

References

Abdous, B., Theodorescu, R. (1992) Note on the spatial quantile of a random vector, *Statistics and Probability Letters* **13**: 333–336.

Breckling, J., Kokic, P., Lübke, O. (2001) A note on multivariate M -quantiles, *Statistics and Probability Letters* **55**: 39–44.

Cai, J.-J., Einmahl, J.H.J., de Haan, L. (2011) Estimation of extreme risk regions under multivariate regular variation, *Annals of Statistics* **39**(3): 1803–1826.

Chaouch, M., Goga, C. (2010) Design-based estimation for geometric quantiles with application to outlier detection, *Computational Statistics and Data Analysis* **54**: 2214–2229.

Chaudhuri, P. (1996) On a geometric notion of quantiles for multivariate data, *Journal of the American Statistical Association* **91**(434): 862–872.

Cheng, Y., De Gooijer, J.G. (2007). On the u -th geometric conditional quantile. *J. Statist. Plann. Inference* **137**, 1914–1930.

S. Girard and G. Stupfler (2014) Asymptotic behaviour of extreme geometric quantiles and their estimation under moment conditions. Available at <http://hal.inria.fr/hal-01060985>.

S. Girard and G. Stupfler (2014) Extreme geometric quantiles in a multivariate regular variation framework, working paper.

Koltchinskii, V.I. (1997) M-estimation, convexity and quantiles, *Annals of Statistics* **25**(2): 435–477.

Serfling, R. (2002) Quantile functions for multivariate analysis: approaches and applications, *Statistica Neerlandica* **56**(2): 214–232.

Zuo, Y., Serfling, R. (2000) General notions of statistical depth function, *Annals of Statistics* **28**: 461–482.