

THESE de DOCTORAT de l'UNIVERSITE PARIS 6

Spécialité : Probabilités.

Présentée par Pierre COHORT
pour obtenir le grade de docteur de l'université PARIS 6.

Titre : "Sur quelques problèmes de quantification."

soutenue le 25 Janvier 2000

devant le jury composé de :

- Vlad BALLY (examineur)
- Jean Claude FORT (examineur)
- Jean JACOD (examineur)
- Harald LUSCHGY (rapporteur)
- Gilles PAGÈS (directeur)
- Denis TALAY (rapporteur)

Introduction.

In some applications, such as source compression, it is useful to quantize a real signal X with N values $y_1 < \dots < y_N$. To this end, Voronoï quantization is generally used : X is approximated by y_i when i is the smallest index which checks $|y_i - X| = \min_{1 \leq j \leq N} |y_j - X|$. The error induced by this approximation is usually modeled by $\min_{1 \leq j \leq N} V(y_j - X)$ where the loss function V is even, zero at 0, strictly increasing on \mathbf{R}_+ and the quality of the quantizer (y_1, \dots, y_N) is measured by the distortion, namely, the average error

$$D_{\mu, N, V}^{y_1, \dots, y_N} := \mathbf{E} \left(\min_{1 \leq j \leq N} V(y_j - X) \right)$$

induced on the data.

Useful quantizers are those that induce a small distortion and especially a distortion equal to $\inf D_{\mu, N, V}^{y_1, \dots, y_N}$. The existence of a quantizer x^* satisfying $D_{\mu, N, V}(x^*) = \inf D_{\mu, N, V}^{y_1, \dots, y_N}$ is a well known fact relying on Fatou's Lemma. Such a quantizer is not in general explicitly known and has to be numerically computed. Unfortunately, classical algorithms optimizing $D_{\mu, N, V}$ (for instance, Lloyd I, Lloyd-Max, stochastic gradient, see [5], [10] and [7]) can get trapped into sub-optimal stationary quantizers, *i.e.*, some zeros of $\nabla D_{\mu, N, V}$ which do not lie in $\text{argmin} D_{\mu, N, V}$. For this reason, it has been natural to look for couples (V, μ) ensuring the uniqueness of a stationary quantizer (which implies of course the uniqueness of a locally optimal quantizer).

The earliest result in that direction was obtained by Fleisher (see [2]) in 1964. The author proved the uniqueness for the quadratic distortion (*i.e.* $V(x) = (x)^2$) and an absolutely continuous law having a strict log-concave differentiable density. In 1982, Trushkin [10] proposed some more general conditions on (V, μ) ensuring uniqueness. For instance, he allows V to depend both on the signal and on the error quantization. Trushkin showed that couples $((\cdot)^2, \mu)$ and $(|\cdot|, \mu)$ where μ has a log-concave density satisfy his conditions and then ensure uniqueness. At the same time, Kieffer [5], extended Trushkin's result about the quadratic loss to more general loss functions, namely, V convex, continuously differentiable and even. As far as we know, the last result has been done by Trushkin. In [11], the author proves the uniqueness for log-concave densities and convex even loss functions, removing then the differentiability assumption made by Kieffer.

The methods used by Trushkin and Kieffer are constructive in the sense that log-concavity ensures convergence of some algorithm toward the unique stationary quantizer and then yields a procedure to numerically compute it ; in Trushkin's method, log-concavity ensures monotonicity of the Lloyd-max algorithm (see [11]), which proves uniqueness and yields a procedure that makes the algorithm converge. In Kieffer's method, log-concavity grants that the map defining the Lloyd I algorithm is a contraction. Then it ensures uniqueness via a fixed point principle, together with convergence of the Lloyd I algorithm toward the unique stationary quantizer.

More recently and independently, Lamberton and Pagès [7] proposed a geometrical approach to the uniqueness problem. The initial goal was to prove the convergence in the Kushner & Clark sense of the Kohonen algorithm (Bouton & Pagès, [1], Fort & Pagès, [3]). This convergence is obtained via the existence of a stable attracting regular basin for the equilibria of the mean function h of the algorithm. Such existence has been proved under strict log-concavity of the density function f by inspecting the spectrum of $\nabla h(x)$ for all $x \in \{h = 0\}$ (see [1] for the two neighbour Kohonen algorithm and [3] for the Kohonen

algorithm with general neighborhood function). In the case of the 0-neighbour Kohonen algorithm, it turns out that $h = \nabla D_{\mu, N, \|\cdot\|^2}$ and the attractivity of equilibria implies their strict minimality. This fact suggested to the authors a method to prove uniqueness of a stationary quantizer, based on the Mountain Pass Lemma. Indeed, the Mountain Pass Lemma says that a function $L : \mathbf{R}^d \rightarrow \mathbf{R}$, coercive, continuously differentiable, having two strict minima necessarily admits a third critical point which is not locally optimal. So, the authors used strict minimality of all stationary quantizer together with a version of the Mountain Pass Lemma to obtain uniqueness. Then, the log-concavity ensures uniqueness and the almost sure convergence (when $\text{supp}(f)$ is compact) of the 0-neighbour Kohonen algorithm toward the unique stationary quantizer.

Pagès & Fort's result is less general than the Trushkin's one since it deals only with quadratic distortion and strictly log-concave density function. So, it seems natural to extend the geometrical method to more general loss function, at least embodying Trushkin's ones and to deal with the well known critical case of some piecewise log-affine density functions.

Next, we give the outline of the paper. The definitions of conditions (A), (C), (L) and (M) are provided in the next section.

Our method consists in

- (a) showing the strict minimality of all critical points of $D_{f, N, V}$,
- (b) applying the Mountain Pass Lemma : if $\nabla D_{f, N, V}^{y^1} = 0$ and if $\nabla D_{f, N, V}^{y^2} = 0$, step (a) shows strict minimality of y^1 et y^2 and the Mountain Pass Lemma then ensures the existence of $y^3 \in \{\nabla D_{f, N, V} = 0\}$ which is not locally optimal, which contradicts (a).

We will show that $D_{f, N, V}$ is continuously differentiable provided that (V) and (C) hold (Proposition 0.2, section 0.7). Then, assuming in addition Assumption (M), we will derive a twice partial differentiability property for $D_{f, N, V}$ (Proposition 0.3, section 0.7). At this stage, the proof of uniqueness splits in two parts :

First, we will assume that f satisfies (L) but not (A). In that case, we will analyse the spectrum of some partial Hessian on the set $\{\nabla D_{f, N, V} = 0\}$ to show that critical points of the distortion are not degenerate. Therefore, a second order criterion will permit us to derive step (a) (Proposition 0.4, section 0.8). Next, step (b) will follow from the Mountain Pass Lemma (Proposition 0.5, section 0.9), completing then the proof of uniqueness of a locally optimal quantizer under the assumptions (V), (C) and $(L) \setminus (A)$.

Second, we will assume that f satisfies (A). In that case, it turns out that critical points of the distortion may be degenerate so that the second order criterion fails. We will then use an approximation procedure to show step (a). As in the first part, step (b) will still follow from the Mountain Pass Lemma, (Proposition 0.6, section 0.10), completing then the proof of uniqueness of a locally optimal quantizer under the assumptions (V), (C) and (A). From Proposition 0.5 and Proposition 0.6, we will then obtain uniqueness under Assumption (V), (C) and (L) Theorem 0.11.

In the case that the unique stationary quantizer of $D_{f, N, V}$ is a degenerate critical point, it is interesting to provide some information about the behaviour of $D_{f, N, V}$ near its absolute minimum. To this end, we will carry out a third order expansion of $D_{f, N, V}$ (Proposition 0.7, section 0.11).

Finally, we will give a simple procedure to generate counter-examples to uniqueness of a locally optimal quantizer (section 0.12).

0.6 Notations and definitions.

Let m, M be two real numbers such that $-\infty \leq m < M \leq +\infty$ and let O_N be the set

$$O_N := \{y \in \mathbf{R}^N; m \leq y_1 < \dots < y_N \leq M\}$$

If $y \in O_N$, one define the i^{th} Voronoï cell $C_i(y)$ of $y \in O_N$, by

$$C_i(y) := \{u \in \mathbf{R}; |y_i - u| \leq |y_j - u|, \forall j \neq i\} = [\tilde{y}_i, \tilde{y}_{i+1}]$$

where $\tilde{y}_1 := m$, $\tilde{y}_i = \frac{y_{i-1} + y_i}{2}$ $i \in \{2, \dots, N\}$ and $\tilde{y}_{N+1} := M$.

The Voronoï quantization of a random variable X by the quantizer $y \in O_N$ is defined by

$$Q_N(X) := \sum_{i=1}^N y_i 1_{C_i(y)}(X) \quad .$$

Quantizing X by $Q_N(X)$ produces a quantization error, designed here by

$$V(Q_N(X) - X)$$

where $V: \mathbf{R} \mapsto \mathbf{R}_+$ is a loss function satisfying

$$(\mathcal{V}) \equiv \begin{cases} (i) & V(u) = \int_0^u v(s) ds \\ (ii) & v \text{ is càdlàg, odd, bounded on compact sets} \\ (iii) & \forall z_0 \in \mathbf{R}^*, v_{\pm}(z_0) := \lim_{z \rightarrow z_0^{\pm}} v(z) > 0 \end{cases} .$$

In section 0.12.2 We give some examples of loss functions checking Assumption \mathcal{V} .

Assumption (\mathcal{V}) implies that $V(0) = 0$, V is even, strictly increasing on \mathbf{R}^+ , right and left differentiable with $V'^{\pm}(z) = v_{\pm}(z)$.

The set

$$Disc(v) := \{z \in \mathbf{R}; v \text{ discontinuous at } z\}$$

is countable and V is differentiable on ${}^c Disc(v)$.

Let μ be the distribution of the random variable X . We will assume that

$$(\mathcal{M}) \equiv \begin{cases} (i) & \mu = f\lambda \text{ with absolutely continuous } f \\ (ii) & \text{supp}(\mu) = [m, M] \end{cases} .$$

The mean quantization error produced by a quantizer $y \in O_N$ is then

$$\begin{aligned} D_{f, N, V}^{y_1, \dots, y_N} &:= \mathbf{E}(V(Q_N(X) - X)) \\ &= \int_m^M \min_{1 \leq i \leq N} V(y_i - u) f(u) du \\ &= \sum_{i=1}^N \int_{\tilde{y}_i}^{\tilde{y}_{i+1}} V(y_i - u) f(u) du \end{aligned}$$

$D_{f,N,V}$ is called the distortion function. Note that if we define $D_{f,N,V}$ by the last equality, the density f need not to be a probability density. The quantizer x may then be viewed as a quantizer of the measure μ .

In some cases, we will not suppose $\mu = f\lambda$; the distortion will then be written $D_{\mu,N,V}$.

In order to ensure that the distortion is finite and sufficiently smooth, it is necessary to assume that

$$(C) \equiv \begin{cases} (i) & \forall x \in \mathbf{R}, \int_m^M V(x-u) \mu(du) < +\infty \\ (ii) & \forall \text{ compact } K, K \subset \mathbf{R}, u \mapsto \sup_{x \in K} |v(x-u)| \in L^1(\mu) \\ (iii) & \forall \text{ compact } K, K \subset \mathbf{R}, u \mapsto \sup_{x \in K} |v(x-u)| |f'(u)| \in L^1(du) \end{cases}$$

Note That $C(ii)$ is well defined even if μ does not check \mathcal{M} .

The main assumption on μ to get uniqueness will be

$$(\mathcal{L}) \equiv \log(f) \text{ is concave on }]m, M[.$$

The density functions that check (\mathcal{L}) are said to be log-concave. The fact that $(\mathcal{L}) \implies (\mathcal{M})$ is straightforward.

In our method, solving the uniqueness problem will require to treat separately the case of log-affine density functions having a finite number of log-affinity intervals, that is, the density functions satisfying

$$(\mathcal{A}) \equiv \begin{cases} f \text{ is continuous.} \\ \exists \theta_2 < \dots < \theta_k \text{ such that } f(u) := \sum_{i=1}^k e^{a_i + b_i u} \mathbf{1}_{] \theta_i, \theta_{i+1}[}(u) \\ \text{where } \theta_1 = -\infty \text{ and } \theta_{k+1} = +\infty \end{cases}$$

One has $(\mathcal{A}) \implies (\mathcal{L})$. The notation $(\mathcal{L}) \setminus (\mathcal{A})$ will denotes the density functions that satisfy (\mathcal{L}) but not (\mathcal{A}) . One can verify that $(\mathcal{L}) \setminus (\mathcal{A})$ holds if and only if one of the following conditions holds

- (a) f is strictly log-concave on an open set $\mathcal{O} \subset]m, M[$.
- (b) f is piecewise log-affine with an infinite number of log-affinity interval.
- (c) $f(m) + f(M) > 0$.

In the sequel, a measure μ on $(\mathbf{R}, \mathcal{B}(\mathbf{R}))$ is said to be continuous if $\mu(\{x\}) = 0$ for every $x \in \mathbf{R}$.

0.7 Regularity of the distortion.

Proposition 0.2. Assume that μ is continuous and that Assumptions (V) et (C) are satisfied. Then, $D_{\mu,N,V}$ is continuously differentiable on O_N and

$$\nabla D_{\mu,N,V}^{y_1, \dots, y_N} = \left(\int_{\bar{y}_i}^{\bar{y}_{i+1}} v(y_i - u) \mu(du) \right)_{1 \leq i \leq N}$$

Proof. Let $\mathbf{y} \in O_N$; since μ is continuous, the set $\text{Disc}(v)$ and the boundaries of the Voronoï cells of \mathbf{y} have zero mass. Assumption (V) then yields the $\mu(du)$ -almost sure differentiability of $\mathbf{z} \mapsto \min_{1 \leq i \leq N} V(z_i - u)$ at \mathbf{y} with gradient $(1_{C_i(\mathbf{y})} v(y_i - u))_{1 \leq i \leq N}$.

But, for $\|\mathbf{y} - \mathbf{z}\| \leq 1$,

$$\begin{aligned} \left| \min_{1 \leq i \leq N} V(z_i - u) - \min_{1 \leq i \leq N} V(y_i - u) \right| &\leq \max_{1 \leq i \leq N} |V(z_i - u) - V(y_i - u)| \\ &\leq \max_{1 \leq i \leq N} \left(\left(\sup_{c \in [y_i - 1, y_i + 1]} |v(c - u)| \right) |y_i - z_i| \right). \end{aligned}$$

Assumption (C) (ii) yields that $u \mapsto \sup_{c \in [y_i - 1, y_i + 1]} |v(c - u)| \in L^1(\mu)$. The Lebesgue dominated convergence Theorem finally yields the differentiability of $D_{\mu, N, V}$ with gradient

$$\nabla D_{\mu, N, V}^{y_1, \dots, y_N} = \left(\int_{\tilde{y}_i}^{\tilde{y}_{i+1}} v(y_i - u) \mu(du) \right)_{1 \leq i \leq N}.$$

Then, Assumption (C) (ii) ensures the dominated convergence of $v(c_m - u)$ toward $v(c - u)$ when $c_m \rightarrow c$. Hence, $\nabla D_{f, N, V}$ is continuous on O_N . \diamond

Now, assume that μ satisfies (M).

The set Ω_i will denote

$$\Omega_i := \{\mathbf{y} \in O_N; y_j - \tilde{y}_j \notin \text{Disc}(v); \quad \forall j \in \{i, i+1\}\}$$

and the set Ω will denote $\bigcap_{i=1}^N \Omega_i$.

For $\mathbf{y} \in O_N$ and $\alpha \in \{+1, -1\}^N$, we consider the open (affine) cone

$$C_\alpha(\mathbf{y}) := \left\{ \begin{array}{l} z \in O_N; \quad \text{sign}(z_1 - y_1) = \alpha_1 \\ \text{sign}((z_i - y_i) - (z_{i-1} - y_{i-1})) = \alpha_i, \quad 2 \leq i \leq N-1 \\ \text{sign}(z_N - y_N) = \alpha_N \end{array} \right\}$$

where $\text{sign}(x) := 1_{\mathbf{R}_+^*}(x) - 1_{\mathbf{R}_-^*}(x)$.

Definition 0.3. Let $G : O_N \mapsto \mathbf{R}$, $\mathbf{y} \in O_N$ and $\alpha \in \{+, -\}^N$. We say that G is differentiable relatively to $C_\alpha(\mathbf{y})$ if G is differentiable at \mathbf{y} for the topology induced by \mathbf{R}^N on $C_\alpha(\mathbf{y})$.

We have the

Proposition 0.3. Let $\mathbf{y} \in O_N$ and let $\alpha \in \{+1, -1\}^N$. Assume that the assumptions (V), (C) and (M) hold. Then, the function $D_{f, N, V}$ is twice differentiable at \mathbf{y} relatively to $C_\alpha(\mathbf{y})$ with Hessian

$$\nabla_\alpha^2 D_{f, N, V}^{y_1, \dots, y_N} = \begin{pmatrix} a_1 + b_2 + \xi_1 & -b_2 & 0 & \cdots & 0 & \cdots & 0 \\ -b_2 & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 & \ddots & 0 \\ \vdots & \ddots & -b_i & a_i + b_i + b_{i+1} & -b_{i+1} & \ddots & \vdots \\ 0 & \ddots & 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & -b_N \\ 0 & \cdots & 0 & \cdots & 0 & -b_N & a_N + b_N - \xi_N \end{pmatrix}.$$

$$\text{where } \begin{cases} a_i = \int_{\tilde{y}_i}^{\tilde{y}_{i+1}} v(y_i - u) f'(u) du & 1 \leq i \leq N \\ b_i = \frac{1}{2} v_{\alpha_i}(y_i - \tilde{y}_i) f(\tilde{y}_i) & 2 \leq i \leq N \\ \xi_1 = \mathbf{1}_{\mathbf{R}}(m) v_{\alpha_1}(y_1 - m) f(m) \\ \xi_N = \mathbf{1}_{\mathbf{R}}(M) v_{\alpha_N}(y_N - M) f(M) \end{cases}$$

Proof. Remember that

$$\Omega_j := \{y \in O_N; y_j - \tilde{y}_j \notin \text{Disc}(v); \quad \forall j \in \{i, i+1\}\}$$

Let $i \in \{1, \dots, N\}$ and let $z \in \Omega_i \cap O_N$. The partial derivative $d_{i,j}^2 D_{f,N,V}^{z_1, \dots, z_N}$ exists and vanishes for $j > i+1$ and $j < i-1$. The continuity of the function v at $z_i - \tilde{z}_i$ yields the existence of $d_{i,i-1}^2 D_{f,N,V}^{z_1, \dots, z_N}$. Indeed,

$$\begin{aligned} d_{i,i-1}^2 D_{f,N,V}^{z_1, \dots, z_N} &= d_{i-1} \left(\int_{\tilde{z}_i}^{\tilde{z}_{i+1}} v(z_i - u) f(u) du \right) \\ &= -\frac{1}{2} v(z_i - \tilde{z}_i) f(\tilde{z}_i) . \end{aligned}$$

Similarly, we obtain the existence of $d_{i,i+1}^2 D_{f,N,V}^{z_1, \dots, z_N}$ for $1 \leq i \leq N-1$ with

$$d_{i,i+1}^2 D_{f,N,V}^{z_1, \dots, z_N} = \frac{1}{2} v(z_i - \tilde{z}_{i+1}) f(\tilde{z}_{i+1}) = -\frac{1}{2} v(z_{i+1} - \tilde{z}_{i+1}) f(\tilde{z}_{i+1}) .$$

Besides, we have

$$\begin{aligned} d_{i,i}^2 D_{f,N,V}^{z_1, \dots, z_N} &= d_i \left(\int_{\tilde{z}_i}^{\tilde{z}_{i+1}} v(z_i - u) f(u) du \right) \\ &= d_i \left(\int_{z_i - \tilde{z}_{i+1}}^{z_i - \tilde{z}_i} v(u) f(z_i - u) du \right) . \end{aligned}$$

f being absolutely continuous, Assumption C (ii), C (iii) and the Lebesgue dominated convergence theorem ensures the existence of $d_{i,i}^2 D_{f,N,V}^{z_1, \dots, z_N}$; namely

$$d_{i,i}^2 D_{f,N,V}^{z_1, \dots, z_N} = \begin{cases} \frac{1}{2} v(z_i - \tilde{z}_i) f(\tilde{z}_i) + \frac{1}{2} v(z_{i+1} - \tilde{z}_{i+1}) f(\tilde{z}_{i+1}) + \int_{\tilde{z}_i}^{\tilde{z}_{i+1}} v(z_i - u) f'(u) du, \\ \quad \text{for } 2 \leq i \leq N-1. \\ \frac{1}{2} v(z_1 - m) f(m) \mathbf{1}_{\mathbf{R}}(m) + \frac{1}{2} v(z_2 - \tilde{z}_2) f(\tilde{z}_2) + \int_m^{\tilde{z}_2} v(z_1 - u) f'(u) du, \\ \quad \text{for } i = 1. \\ \frac{1}{2} v(z_N - \tilde{z}_N) f(\tilde{z}_N) - \frac{1}{2} v(z_N - M) f(M) \mathbf{1}_{\mathbf{R}}(M) + \int_{\tilde{z}_N}^M v(z_N - u) f'(u) du, \\ \quad \text{for } i = n. \end{cases}$$

Let $\varepsilon > 0$, $\eta > 0$, $y \in \Omega \cap O_N$ and $x \in O_N$ such that $\|y - x\| \leq \eta$. For all $i \in \{1, \dots, N\}$ there exists a countable set $\tau_i \subset [y_i, y_i]$ such that for all $t \in [y_i, y_i] \setminus \tau_i$, $(y_1, \dots, y_{i-1}, t, y_{i+1}, \dots, x_N) \in \Omega_i$. So, the partial derivative $d_i \nabla D_{f,N,V}^{y_1, \dots, y_{i-1}, t, y_{i+1}, \dots, y_N}$ exists for all $t \in [y_i, y_i] \setminus \tau_i$ and, by the bounded increments formula, we obtain

$$\left\| \nabla D_{f,N,V}^{x_1, \dots, x_N} - \nabla D_{f,N,V}^{y_1, \dots, y_N} - \sum_{i=1}^N d_i \nabla D_{f,N,V}^{y_1, \dots, y_N} (y_i - x_i) \right\|$$

$$\leq \sum_{i=1}^N \sup_{t \in [y_i, y_i] \setminus \tau_i} \left\| d_i \nabla D_{f,N,V}^{y_1, \dots, y_{i-1}, t, y_{i+1}, \dots, x_N} - d_i \nabla D_{f,N,V}^{y_1, \dots, y_{i-1}, y_i, y_{i+1}, \dots, x_N} \right\| |y_i - y_i|$$

But, since $\mathbf{y} \in \Omega \cap O_N$, the function $t \mapsto d_i \nabla D_{f,N,V}^{y_1, \dots, y_{i-1}, t, y_{i+1}, \dots, x_N}$ defined on $[y_i, y_i] \setminus \tau_i$ is continuous at $t = y_i$ and then, for sufficiently small η , we have

$$\left\| \nabla D_{f,N,V}^{x_1, \dots, x_N} - \nabla D_{f,N,V}^{y_1, \dots, y_N} - \sum_{i=1}^N d_i \nabla D_{f,N,V}^{y_1, \dots, y_N} (y_i - y_i) \right\| \leq \varepsilon \|\mathbf{y} - \mathbf{x}\|$$

It follows that $D_{f,N,V}$ is twice differentiable at \mathbf{y} , with Hessian $\nabla^2 D_{f,N,V}^{y_1, \dots, y_N}$.

Now, let $\varepsilon, \eta > 0$, $\mathbf{y} \in {}^c\Omega \cap O_N$, $\alpha \in \{+, -\}^N$, $\mathbf{x} \in C_\alpha(\mathbf{y})$ and $\hat{\mathbf{y}} \in \Omega \cap C_\alpha(\mathbf{y})$ such that $\|\mathbf{y} - \hat{\mathbf{y}}\| < \eta$. One has

$$\begin{aligned} \left\| \nabla D_{f,N,V}^{x_1, \dots, x_N} - \nabla D_{f,N,V}^{y_1, \dots, y_N} - \nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N} (\mathbf{x} - \mathbf{y}) \right\| \leq \\ \left\| \nabla D_{f,N,V}^{x_1, \dots, x_N} - \nabla D_{f,N,V}^{\hat{y}_1, \dots, \hat{y}_N} - \nabla^2 D_{f,N,V}^{\hat{y}_1, \dots, \hat{y}_N} (\mathbf{x} - \hat{\mathbf{y}}) \right\| \\ + \left\| \nabla D_{f,N,V}^{y_1, \dots, y_N} - \nabla D_{f,N,V}^{\hat{y}_1, \dots, \hat{y}_N} \right\| + \left\| \nabla^2 D_{f,N,V}^{\hat{y}_1, \dots, \hat{y}_N} \right\| \|\mathbf{x} - \hat{\mathbf{y}}\| \\ + \left\| \nabla^2 D_{f,N,V}^{\hat{y}_1, \dots, \hat{y}_N} - \nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N} \right\| \|\mathbf{x} - \mathbf{y}\| \end{aligned}$$

But, we have

$$\lim_{\substack{\hat{\mathbf{y}} \rightarrow \mathbf{y} \\ \hat{\mathbf{y}} \in C_\alpha(\mathbf{y})}} \nabla^2 D_{f,N,V}^{\hat{y}_1, \dots, \hat{y}_N} = \nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N}$$

In addition, we know that $D_{f,N,V}$ is twice differentiable at $\hat{\mathbf{y}} \in \Omega \cap O_N$. So, for sufficiently small η , one obtains

$$\left\| \nabla D_{f,N,V}^{x_1, \dots, x_N} - \nabla D_{f,N,V}^{y_1, \dots, y_N} - \nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N} (\mathbf{x} - \mathbf{y}) \right\| \leq \varepsilon \|\mathbf{x} - \mathbf{y}\|$$

then, $D_{f,N,V}$ is twice differentiable relatively to $C_\alpha(\mathbf{y})$ with Hessian $\nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N}$. \diamond

0.8 Strict minimality of critical points: the non-degenerate case.

In this section, we establish the strict minimality of all critical points for density functions satisfying Assumption $(\mathcal{L}) \setminus (\mathcal{A})$. Piecewise log-affine density function will be noted

$$f(z) = \sum_{i \in J} e^{a_i + b_i z} \mathbf{1}_{[\theta_i, \theta_{i+1}]}(z)$$

where $([\theta_i, \theta_{i+1}])_{i \in J}$ denotes the set of log-affinity intervals of f and J is an interval of \mathbb{Z} . The sequence $(b_i)_{i \in J}$ is necessarily strictly decreasing. The spectrum of a matrix $A \in \mathcal{M}_N(\mathbb{R})$ will be noted $Sp(A)$. When all eigenvalues of A are positive real numbers, we write $A > 0$.

Proposition 0.4. *Let $\alpha \in \{+1, -1\}^N$ and $\mathbf{y} \in O_N$ a critical point of $D_{f,N,V}$. Assume that (\mathcal{V}) and (\mathcal{C}) hold. Assume that $(\mathcal{L}) \setminus (\mathcal{A})$ holds, that is, one of the following assumptions holds*

- (a) f is strictly log-concave on an open set $\mathcal{O} \subset]m, M[$.
- (b) f is piecewise log-affine and J is not finite.
- (c) $f(m) + f(M) > 0$.

Then, x is a strict minimum of $D_{f,N,V}$.

Proof. Thanks to Proposition 0.3, it is sufficient to prove that

$$\forall \alpha \in \{+, -\}^N, \quad \nabla_\alpha^2 D_{f,N,V}(x) > 0.$$

To this end, we use the classical Gershgorin-Hadamard Lemma (see for instance [13]):

Lemma 0.3 (Gershgorin-Hadamard). *Let A be a real matrix satisfying*

- (α) $\forall i \neq j, a_{ij} \leq 0$
- (β) $\forall i, \sum_j a_{ij} \geq 0$
- (γ) $\forall i, a_{i,i \pm 1} < 0$ ($1 \leq i-1$ et $i+1 \leq N$)
- (δ) $\exists i_0$ such that $\sum_j a_{i_0 j} > 0$

then, the real parts of the eigenvalues of A are strictly positive.

By (\mathcal{V}), we have $v_\pm > 0$ on \mathbf{R}_+^* ; the matrix $\nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N}$ then satisfies (α) and (γ). Next we derive that $\nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N}$ checks (β).

Let $i \in \{2, \dots, N\}$ and let $S_i^\alpha := \sum_{j=1}^N [\nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N}]$. One has

$$S_i^\alpha := \int_{\tilde{y}_i}^{\tilde{y}_{i+1}} v(y_i - u) f'(u) du$$

(in fact, only S_1^α and S_N^α may depend on α). We use Lagrange's method to prove that $S_i^\alpha \geq 0$.

Let a, b, c such that $m \leq a \leq c \leq b \leq M$; one defines the function ϕ by

$$\phi(a, b, c) := \int_a^b v(c-u) f'(u) du - \frac{f'^+(c)}{f(c)} \int_a^b v(c-u) f(u) du.$$

The function ϕ has right partial derivatives with respect to a and b given by

$$\frac{\partial \phi}{\partial a^+}(a, b, c) = v_-(c-a) \left(\frac{f'^+(c)}{f(c)} f(a) - f'^+(a) \right) \quad (17)$$

$$\frac{\partial \phi}{\partial b^+}(a, b, c) = v_-(c-b) \left(f'^+(b) - \frac{f'^+(c)}{f(c)} f(b) \right). \quad (18)$$

Assumption (\mathcal{V}) yields $v_-(c-a) > 0$ and $v_-(c-b) < 0$. The log-concavity of f (Assumption (\mathcal{L})) then implies that $\frac{f'^+}{f}$ is decreasing; so, the righthand term in (17) (resp. in (18)) is ≤ 0 (resp. ≥ 0).

Finally, $\phi(y_i, y_i, y_i) = 0$ entails that $\phi(\tilde{y}_i, \tilde{y}_{i+1}, y_i) \geq 0$.
 Since $x \in \{\nabla D_{f,N,V} = 0\}$, one has

$$\int_{\tilde{y}_i}^{\tilde{y}_{i+1}} v(y_i - u) f(u) du = 0$$

thus, $S_i^\alpha \geq 0$ and (β) holds. Next we derive that $\nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N}$ checks (δ) .

Assume (a); There is some index i_0 such that $\mathcal{O} \cap]\tilde{y}_{i_0}, \tilde{y}_{i_0+1}[\neq \emptyset$. The function $\frac{f'}{f}$ is strictly decreasing on $\mathcal{O} \cap]\tilde{y}_{i_0}, \tilde{y}_{i_0+1}[$ and then,

$$\begin{aligned} \forall a \in \mathcal{O} \cap]\tilde{y}_{i_0}, y_{i_0}[\quad \frac{\partial \phi}{\partial a^+}(a, \tilde{y}_{i_0+1}, y_{i_0}) &> 0 \\ \forall b \in \mathcal{O} \cap]y_{i_0}, \tilde{y}_{i_0+1}[\quad \frac{\partial \phi}{\partial b^+}(\tilde{y}_{i_0}, b, y_{i_0}) &> 0 \end{aligned}$$

So, we conclude that $S_{i_0}^\alpha > 0$.

When (b) is satisfied, we prove in the same way that $S_{i_0}^\alpha > 0$ using the inequality $b_i > b_{i+1}$.
 At last, Assumption (c) implies that $S_1^\alpha > 0$ or $S_N^\alpha > 0$.

For the three cases, (δ) is satisfied and then, Lemma 0.3 applies and yields

$$\text{Re} \left(\text{Sp} \left(\nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N} \right) \right) > 0$$

The matrix $\nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N}$ being symmetric, we have $\text{Sp} \left(\nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N} \right) \subset \mathbf{R}$. This completes the proof. \diamond

0.9 Uniqueness : the non-degenerate case.

The goal of this section is to prove and then apply a version of the Mountain Pass Lemma adapted for the distortion setting. The classical Mountain Pass Lemma reads as follows

Theorem 0.9. *If L is continuously differentiable, if L is coercive, i.e. $L(x) \rightarrow +\infty$ as $\|x\| \rightarrow +\infty$, and if $\nabla D_{f,N,V}$ has two zeros which are strict minima then $\nabla D_{f,N,V}$ has also a third zero which is not a local minimum.*

Proof. See [10] pp. 66-68. \diamond

In our setting, the distortion $D_{f,N,V}$ is not defined on \mathbf{R}^N but only on O_N . So, we need a version of theorem 0.9 that can be applied here. We propose the

Theorem 0.10. *Let K be a compact and convex set of \mathbf{R}^N with non-empty interior O . If $L : K \rightarrow \mathbf{R}$ is continuously differentiable on O , if ∇L has a continuous extension on K (still denoted ∇L) and if, for sufficiently small $\varepsilon > 0$, $(I_d - \varepsilon \nabla L)(O) \subset O$, then, the conclusion of the Theorem 0.9 holds.*

Proof. A proof is provided in the annex. \diamond

Next proposition proves the uniqueness under the assumptions of Proposition 0.4.

Proposition 0.5. *Assume that (V) and (C) hold. Assume that $(\mathcal{L}) \setminus (\mathcal{A})$ holds, that is, that one of the following assumptions holds*

- (a) f is strictly log-concave on an open set $\mathcal{O} \subset]m, M[$,
- (b) f is piecewise log-affine and J is not finite,
- (c) $f(m) + f(M) > 0$.

Then, the distortion $D_{f,N,V}$ has a unique critical point (in particular, we get the uniqueness of a locally optimal quantizer).

Proof. As a first step, we will assume that $[m, M]$ is compact in order to apply Theorem 0.10. The extension to non-compactly supported distributions will rely on an approximation procedure (step 2).

Step 1 : The set

$$\overline{O}_N = \{m \leq y_1 \leq \dots \leq y_N \leq M\}$$

is compact and connected with non-empty interior. The coordinates of $\nabla D_{f,N,V}$ are defined and continuous on \overline{O}_N ; the gradient $\nabla D_{f,N,V}$ has then a continuous extension on \overline{O}_N (still denoted $\nabla D_{f,N,V}$).

The function v is bounded on the compact set $[m, M]$ (Assumption (V)) and f is bounded since log-concave. Thus, for $\varepsilon < \frac{1}{\|v\|_\infty \|f\|_\infty}$ and $i \in \{1, \dots, N-1\}$, we have

$$\begin{aligned} y_{i+1} - y_i &> \varepsilon \|v\|_\infty \|f\|_\infty (y_{i+1} - y_i) \\ &> \varepsilon \left(\int_{\tilde{y}_{i+1}}^{y_{i+1}} v(y_{i+1} - u) f(u) du - \int_{y_i}^{\tilde{y}_{i+1}} v(y_i - u) f(u) du \right) \\ &> \varepsilon \left([\nabla D_{f,N,V}^{y_1, \dots, y_N}]_{i+1} - [\nabla D_{f,N,V}^{y_1, \dots, y_N}]_i \right) \end{aligned}$$

and then, $(I_d - \varepsilon \nabla D_{f,N,V})(O_N) \subset O_N$.

We can then apply Theorem 0.10 : if $y^1 \in \{\nabla D_{f,N,V} = 0\}$ and $y^2 \in \{\nabla D_{f,N,V} = 0\}$, y^1 and y^2 are strict minima (Proposition 0.4). So, there is $y^3 \in \{\nabla D_{f,N,V} = 0\}$ which is not a local minima. Now, one notices that $\nabla D_{f,N,V}$ does not vanish on $\partial \overline{O}_N$. Indeed, let $y \in \partial \overline{O}_N$ and let i be the index of the first coordinate of y lying in an aggregate (note that $i < N$); since $y_i = \tilde{y}_{i+1} = y_{i+1}$, we have

$$[\nabla D_{f,N,V}^{y_1, \dots, y_N}]_i = \int_{\tilde{y}_i}^{y_i} v(y_i - u) f(u) du \neq 0$$

and $\nabla D_{f,N,V}$ does not vanish on $\partial \overline{O}_N$.

Then, Proposition 0.4 is contradicted. Hence, $y^3 \in \{\nabla D_{f,N,V} = 0\} \cap O_N$. Uniqueness of a critical point for the distortion follows.

Step 2 : If $[m, M]$ is not a compact set, set $f_k := f \mathbf{1}_{[-k, k]}$ for every $k \geq 1$. The existence of at least two critical points of $D_{f,N,V}$ implies the existence of at least two strict minima for $D_{f,N,V}$ (Proposition 0.4). But, one can easily check that

$$D_{f_k, n, V} \xrightarrow{U_K} D_{f, N, V}$$

where $\xrightarrow{U_K}$ denotes the uniform convergence on compact sets on O_N .

Then, for sufficiently large k , $D_{f_k, n, V}$ would have at least two strict minima, which contradicts step 1. This completes the proof. \diamond

0.10 The degenerate case.

Once again, (V) and (C) hold. Here, we deal with the case that Assumption (A) holds, that is,

$$f(u) = \sum_{i=1}^k e^{a_i + b_i u} \mathbf{1}_{[\theta_i, \theta_{i+1}[}(u) \quad \text{with } \theta_1 = m = -\infty, \quad \theta_{k+1} = M = +\infty$$

where k is the number of log-affinity intervals of the density function f .

For these density functions, two kinds of critical points may appear :

- Type I : $\exists j \in \{2, \dots, k\}, \theta_j \notin \{\tilde{y}_i, i = 2, \dots, N\}$

In this case, $\frac{f'}{f}$ is strictly decreasing on an interval $]\tilde{y}_{i_0}, \tilde{y}_{i_0+1}[$ and

$$\forall \alpha \in \{+, -\}^N, \quad S_{i_0}^\alpha := \sum_{j=1}^N \left[\nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N} \right]_{i_0 j} > 0, \quad ,$$

then, following the proof of Proposition 0.4, x is a strict minimum of the distortion.

- Type II : $\forall j \in \{2, \dots, k\}, \theta_j \in \{\tilde{y}_i, i = 2, \dots, N\}$

This property implies that $\frac{f'}{f}$ is constant on each $]\tilde{y}_i, \tilde{y}_{i+1}[$; we have then

$$\forall \alpha \in \{+, -\}^N, \forall i \in \{1, \dots, N\} \quad S_i^\alpha := \sum_{j=1}^N \left[\nabla_\alpha^2 D_{f,N,V}(x) \right]_{ij} = 0 \quad .$$

So, $\forall \alpha \in \{+, -\}^N, \quad \det \left(\nabla_\alpha^2 D_{f,N,V}^{y_1, \dots, y_N} \right) = 0$; thus, the second order criterion fails.

The simplest example of a critical point of type II is the two level quantizer $(-1, 1)$ for the Laplace law $\frac{1}{2} (\mathbf{1}_{\mathbf{R}_-}(u) e^u + \mathbf{1}_{\mathbf{R}_+}(u) e^{-u}) du$.

Proposition 0.6. Assume that (V), (C) and (A) holds. Then, $D_{f,N,V}$ has a unique critical point (in particular, we get the uniqueness of a locally optimal quantizer).

Proof. It amounts to show that all critical points are strict minima. To this end, it suffices to deal with critical points of type II. Let us notice that there is at most a finite number of such points. Indeed, let x be a critical point of type II and let $i \in \{1, \dots, k\}$. One has

$$\left(\left[\nabla D_{f,N,V}^{y_1, \dots, y_N} \right]_j \right)_{n_{i-1}+1 \leq j \leq N_i} = \nabla D_{f \mathbf{1}_{[\theta_i, \theta_{i+1}[}^{n,V}}(x_{n_{i-1}+1}, \dots, x_{N_i})$$

where $n_i := |\{j; y_j < \theta_{i+1}\}|$ ($n_0 := 0$).

But, thanks to Proposition 0.5, $D_{f \mathbf{1}_{[\theta_i, \theta_{i+1}[}^{n,V}}$ has only one critical point. Then, the distortion has at most $\binom{n}{n_1, \dots, n_k}$ critical points of type II. Consequently, such critical points are isolated in $\{\nabla D_{f,N,V} = 0\}$. Thus, it remains to prove minimality of these quantizers.

Let x be a critical point of type II. For $i \in \{1, \dots, k\}$, let $p_i := \max \{j; x_j < \theta_{i+1}\}$.

We will construct a density function f_ε such that

1. $D_{f_\varepsilon, n, V}$ has a unique critical point (absolute minimum), say x^ε .
2. $x^\varepsilon \rightarrow x$ as $\varepsilon \downarrow 0$.
3. $D_{f_\varepsilon, n, V} \xrightarrow{U_K} D_{f, N, V}$ as $\varepsilon \downarrow 0$.

Let $\varepsilon > 0$. Since for all $i \in \{1, \dots, p_1 - 1\}$, $\int_{\tilde{y}_i}^{\tilde{y}_{i+1}} v(y_i - u) e^{a_1 + b_1 u} du = 0$, one has

$$\int_{\tilde{y}_i - \varepsilon}^{\tilde{y}_{i+1} - \varepsilon} v(y_i - \varepsilon - u) e^{a_1 + b_1 u} du = 0.$$

We set $\delta_i := \varepsilon$ and $\delta_i := \varepsilon$ for $i \in \{1, \dots, p_1\}$ and $\eta_{p_1} := \varepsilon$. In the same way,

$$\int_{\tilde{y}_{p_1} - \varepsilon}^{\tilde{y}_{p_1+1} - \varepsilon} v(x_{p_1} - \varepsilon - u) e^{a_2 + b_2 u} du = 0.$$

But, for all $u \in]\tilde{y}_{p_1} - \varepsilon, \tilde{y}_{p_1}[$, one has $e^{a_1 + b_1 u} < e^{a_2 + b_2 u}$ which implies

$$\int_{\tilde{y}_{p_1} - \varepsilon}^{\tilde{y}_{p_1+1} - \varepsilon} v(x_{p_1} - \varepsilon - u) f(u) du < 0.$$

Thus, there exists $\eta_{p_1+1} > \varepsilon$ such that

$$\int_{\tilde{y}_{p_1} - \varepsilon}^{\tilde{y}_{p_1+1} - \eta_{p_1+1}} v(x_{p_1} - \varepsilon - u) f(u) du = 0.$$

and we define $\delta_{p_1+1} := 2\eta_{p_1+1} - \delta_{p_1} > \eta_{p_1+1}$. By induction, one can define two sequences of positive real numbers $(\eta_i)_{1 \leq i \leq N}$, $(\delta_i)_{1 \leq i \leq N}$ depending only on ε and satisfying

$$(\alpha) \int_{\tilde{y}_i - \eta_i}^{\tilde{y}_{i+1} - \eta_{i+1}} v(y_i - \delta_i - u) f(u) du = 0 \quad \forall i \in \{1, \dots, N-1\},$$

$$(\beta) \delta_{i+1} = 2\eta_i - \delta_i \quad \forall i \in \{1, \dots, N-1\},$$

$$(\gamma) \delta_i > \eta_i \text{ and } \eta_i > \delta_{i-1} \quad \forall i \in \{p_1 + 1, \dots, N\}.$$

Note that for all $i \in \{1, \dots, N\}$, η_i and δ_i are continuous functions of ε .

From (γ) , one has $\delta_N > \eta_N$. So,

$$\int_{\tilde{y}_N - \eta_N}^{+\infty} v(y_N - \delta_N - u) f(u) du < 0$$

Then, there exists $M_\varepsilon \in \mathbb{R}$ such that

$$\int_{\tilde{y}_N - \eta_N}^{M_\varepsilon} v(y_N - \delta_N - u) f(u) du = 0 \quad (19)$$

Now, we define $f_\varepsilon := f \mathbf{1}_{]-\infty, M_\varepsilon]}$ and $x^\varepsilon := (y_i - \delta_i)_{1 \leq i \leq N}$. The density function f_ε is continuous, log-concave and satisfies Assumption (c) of Proposition 0.3. As a result, the distortion $D_{f_\varepsilon, n, V}$ has a unique critical point which is the absolute minimum. But, from (α) , (β) and (19), x^ε is a Voronoï quantizer and $x^\varepsilon \in \{\nabla D_{f_\varepsilon, n, V} = 0\}$. Then, x^ε is the unique critical point, absolute minimum of $D_{f_\varepsilon, n, V}$.

Recall that for all $i \in \{1, \dots, N\}$, η_i and δ_i are continuous functions of ε and that $\delta_i(0) = \eta_i(0) = 0$. Consequently, $x^\varepsilon \rightarrow x$ as $\varepsilon \downarrow 0$ and $M_\varepsilon \rightarrow +\infty$ as $\varepsilon \downarrow 0$. Finally, one can easily show that $D_{f_\varepsilon, n, V} \xrightarrow{U_K} D_{f, N, V}$ as $\varepsilon \downarrow 0$. So, x is a minimum of the distortion $D_{f, N, V}$. As already said, x is isolated in $\{\nabla D_{f, N, V} = 0\}$ which implies that x is a strict minimum of $D_{f, N, V}$. We conclude the proof as in Proposition 0.4. \diamond

To sum up, we have

Theorem 0.11. *Assume that V and μ satisfy (\mathcal{V}) , (\mathcal{C}) and (\mathcal{L}) . Then, $D_{f,N,V}$ has a unique critical point (in particular, we get the uniqueness of a locally optimal quantizer).*

0.11 Order of critical points.

When the unique critical point of the distortion is of type II, , the matrix $\nabla_{\alpha}^2 D_{f,N,V}^{y_1, \dots, y_N}$ is degenerate for all $\alpha \in \{+1, -1\}^N$. We will then deal with the local behaviour of $D_{f,N,V}$ in the neighbourhood of such x .

Let $\mathbf{1} := (1, \dots, 1) \in \mathbb{R}^N$ and recall that $p_j := \min \{l; x_l > \theta_j\}$.

Proposition 0.7. *If x is of type II, one has $D_{f,N,V}^{x+h\mathbf{1}} = D_{f,N,V}^{y_1, \dots, y_N} + |O(h^3)|$. More precisely,*

- for $h > 0$, $D_{f,N,V}^{x+h\mathbf{1}} = D_{f,N,V}^{y_1, \dots, y_N} + K^+ h^3 + o(h^3)$ where

$$K^+ = \sum_{j=2}^k (b_{j-1} - b_j) v_+ (y_{p_j-1} - \tilde{y}_{p_j}) f_{j-1}(\tilde{y}_{p_j}) > 0$$

- for $h < 0$, $D_{f,N,V}^{x+h\mathbf{1}} = D_{f,N,V}^{y_1, \dots, y_N} - K^- h^3 + o(h^3)$ where

$$K^- = \sum_{j=2}^k (b_{j-1} - b_j) v_- (y_{p_j-1} - \tilde{y}_{p_j}) f_{j-1}(\tilde{y}_{p_j}) > 0$$

To prove the proposition, we need the following classical lemma

Lemma 0.4. *Let A be a symmetric tridiagonal matrix with non-zero tridiagonal coefficients, such that $\sum_{j=1}^N a_{ij} = 0$ for $1 \leq i \leq N$. Then, the eigensubspace related to the eigenvalue 0 is one dimensional and spanned by $\mathbf{1}$.*

Proof of the proposition 0.7.

Since $D_{f,N,V}$ is differentiable on O_N , the function $\varphi(h) = D_{f,N,V}(x + h\mathbf{1})$ is differentiable near 0, with derivative $\varphi'(h) = \sum_{i=1}^N \int_{\tilde{y}_i+h}^{\tilde{y}_{i+1}+h} v(y_i + h - u) f(u) du$. The dominated convergence theorem yields the twice differentiability of φ and

$$\varphi''(h) = \sum_{i=1}^N \int_{\tilde{y}_i+h}^{\tilde{y}_{i+1}+h} v(y_i + h - u) f'(u) du.$$

Integrating by parts, we obtain

$$\sum_{i=1}^N \int_{\tilde{y}_i+h}^{\tilde{y}_{i+1}+h} V(y_i+h-u) f'(u) du + \left(\sum_{i=p_{j-1}}^{p_j-2} b_{j-1} \left[\frac{V(y_i-\tilde{y}_i) f_{j-1}(\tilde{y}_i+h) - V(y_i-\tilde{y}_{i+1}) f_{j-1}(\tilde{y}_{i+1}+h)}{V(y_{p_{j-1}}+h-\tilde{y}_{p_j}) f_{j-1}(\tilde{y}_{p_j})} \right] + \sum_{j=2}^{k+1} b_j \left[\frac{V(y_{p_{j-1}}+h-\tilde{y}_{p_j}) f_j(\tilde{y}_{p_j}) - V(y_{p_{j-1}}-\tilde{y}_{p_j}) f_j(\tilde{y}_{p_j}+h)}{V(y_{p_{j-1}}+h-\tilde{y}_{p_j}) f_j(\tilde{y}_{p_j})} \right] \right)$$

which equals

$$\sum_{i=1}^N \int_{\tilde{y}_i+h}^{\tilde{y}_{i+1}+h} V(y_i+h-u) f'(u) du + \left(\sum_{i=p_{j-1}}^{p_j-2} b_{j-1} V(y_i-\tilde{y}_{i+1}) f_{j-1}(\tilde{y}_{i+1}+h) + \sum_{j=2}^{k+1} b_j V(y_{p_{j-1}}+h-\tilde{y}_{p_j}) f_j(\tilde{y}_{p_j}) \right)$$

Another derivation yields

$$\varphi'''(h) = \left(\sum_{i=1}^N \int_{\tilde{y}_i+h}^{\tilde{y}_{i+1}+h} v(y_i+h-u) f'(u) du + \sum_{j=2}^{k+1} \left(\sum_{i=p_{j-1}}^{p_j-2} b_{j-1} v_+(y_i-\tilde{y}_{i+1}) f_{j-1}(\tilde{y}_{i+1}+h) + b_{j-1} v_+(y_{p_{j-1}}+h-\tilde{y}_{p_j}) f_{j-1}(\tilde{y}_{p_j}) + b_j v_+(y_{p_{j-1}}+h-\tilde{y}_{p_j}) f_j(\tilde{y}_{p_j}) \right) \right)$$

When $h = 0$, this term finally equals

$$\sum_{j=2}^k (b_{j-1} - b_j) v_+(y_{p_{j-1}} - \tilde{y}_{p_j}) f_{j-1}(\tilde{y}_{p_j})$$

The sequence $(b_j)_{1 \leq j \leq k}$ is strictly decreasing and v_+ is strictly decreasing on \mathbf{R}_+^* ; we have then $\varphi'''(0) > 0$.

The case $h < 0$ is similar. \diamond

Remark. Proposition 0.7 does not permit to conclude that critical points of type II are strict minima because of the distortion is not, in general, twice continuously differentiable.

0.12 Examples and counter-examples.

0.12.1 Multiplicity of critical points.

In this section, we derive a procedure to generate counter examples to uniqueness of a locally optimal quantizer. Here, the measure μ is only assumed to be continuous and we

extend $D_{\mu,N,V}$ to O_{2n-1} by setting

$$\tilde{D}_{\mu,n,V}(y_1, w_2, x_2, \dots, w_N, y_N) := \sum_{i=1}^N \int_{w_i}^{w_{i+1}} V(y_i - u) \mu(du)$$

with $(y_1, w_2, x_2, \dots, w_N, y_N) \in O_{2n-1}$, $w_1 = -\infty$ and $w_{N+1} = +\infty$. Since V is even and strictly increasing on \mathbf{R}_+ , we notice that

$$D_{\mu,N,V}(y_1, x_2, \dots, y_N) \leq \tilde{D}_{\mu,n,V}(y_1, w_2, x_2, \dots, w_N, y_N)$$

Let $x \in \{\nabla D_{\mu,N,V} = 0\}$. Then, one has

$$\forall i \in \{1, \dots, N\} \quad \int_{\tilde{y}_i}^{\tilde{y}_{i+1}} V(y_i - u) \mu(du) = 0 \quad .$$

So, for all n -uple $(\gamma_1, \dots, \gamma_N) \in \mathbf{R}_+^N$, one obtains

$$\forall i \in \{1, \dots, N\} \quad \int_{\tilde{y}_i}^{\tilde{y}_{i+1}} V(y_i - u) (\gamma_i \mu)(du) = 0$$

that is, x is a critical point of the distortion $D_{\mu_{\gamma,x},n,V}$ where

$$\mu_{\gamma,x} := \sum_{i=1}^N \gamma_i \mu(\cdot \cap C_i(\mathbf{y}))$$

Now, assume that x is an optimal quantizer for μ . We construct the measure $\mu_{\gamma,x}$ as follows :

- If $\mu(C_i(\mathbf{y}))$ is great, choose a small γ_i .
- If $\mu(C_i(\mathbf{y}))$ is small, choose a large γ_i .

This procedure will drastically modify the mass repartition of the measure μ . As a consequence, the levels y_1, \dots, y_N of the μ -optimal quantizer x will probably not match to the mass repartition of $\mu_{\gamma,x}$. So, the quantizer x will probably become sub-optimal for $\mu_{\gamma,x}$. But, as already said, x is still a stationary quantizer for $D_{\mu_{\gamma,x},n,V}$. Finally, $D_{\mu_{\gamma,x},n,V}$ will have at least two stationnary quantizer : x and a $\mu_{\gamma,x}$ -optimal quantizer $x_\gamma^* \neq x$.

We derive this procedure in the following particular case : starting from a stationnary quantizer (not necessarily optimal) of $D_{\mu,N,V}$, we will decrease the mass of only one cell, say, $C_i(\mathbf{y})$. To show that x is not $\mu_{\gamma,x}$ -optimal, we will construct a quantizer y such that

$$D_{\mu_{\gamma,x},n,V}(y) < D_{\mu_{\gamma,x},n,V}(x)$$

by putting the level y_i in another cell, say, $C_j(x)$.

More formally, let $x \in \{\nabla D_{\mu,N,V} = 0\}$ such that $\mu(C_i(\mathbf{y})) > 0 \quad \forall i \in \{1, \dots, N\}$. Let $i, j \in \{2, \dots, N\}$ such that $i < j$.

For $\delta > 0$, set

$$\gamma := (1, \dots, 1, \delta, 1, \dots, 1)$$

where $\gamma_i = \delta$. Define $z = (y_1, w_2, y_2, \dots, w_N, y_N)$ by

- $y_k = x_k$ for $1 \leq k \leq i-1$ et $j+1 \leq k \leq N$.
- $w_k = \tilde{y}_k$ for $2 \leq k \leq i-1$ et $j+1 \leq k \leq N$.
- $w_i = y_i$
- $w_k = \tilde{y}_{k+1}$ for $i+1 \leq k \leq j-1$ et
- $y_k = x_{k+1}$ for $i \leq k \leq j-1$
- $y_j \in]y_j, \tilde{y}_{j+1}[$
- $w_j = \frac{y_{j-1} + y_j}{2}$

We have

$$\begin{aligned} \tilde{D}^{V, \mu_{\gamma, x}}(z) = & \delta \left(\int_{\tilde{y}_i}^{y_i} V(y_{i-1} - u) \mu(du) + \int_{y_i}^{\tilde{y}_{i+1}} V(y_i - u) \mu(du) \right) \\ & - \int_{\tilde{y}_i}^{\tilde{y}_{i+1}} V(y_i - u) \mu(du) - \zeta + D_{\mu_{\gamma, x, n, V}}(x) \end{aligned}$$

where $\zeta > 0$ is the decrease of distortion induced by the introduction of y_j in the set $]y_j, \tilde{y}_{j+1}[$, that is

$$\zeta := \int_{y_j}^{\tilde{y}_{j+1}} V(y_j - u) \mu(du) - \left(\int_{y_{j-1}}^{w_j} V(y_{j-1} - u) \mu(du) + \int_{w_j}^{w_{j+1}} V(y_j - u) \mu(du) \right)$$

We have $\mu(C_i(y)) > 0$; then, for small enough δ ,

$$\tilde{D}_{\mu_{\gamma, x, n, V}}(z) < D_{\mu_{\gamma, x, n, V}}(x)$$

but,

$$D_{\mu_{\gamma, x, n, V}}(y_1, \dots, y_N) \leq \tilde{D}_{\mu_{\gamma, x, n, V}}(z)$$

and then

$$D_{\mu_{\gamma, x, n, V}}(y_1, \dots, y_N) < D_{\mu_{\gamma, x, n, V}}(x).$$

So, the quantizer x is not in the set of optimal quantizer of $D_{\mu_{\gamma, x, n, V}}$ and $D_{\mu_{\gamma, x, n, V}}$ has several critical points.

Example : Let f be the uniform density function on $[0, 1]$ and let $D^{2, f}(y_1, x_2)$ be the quadratic distortion induced by two level quantizers. The unique critical point for $D^{2, f}$ is $(0.25, 0.75)$. Now, let $0 < \delta < 1$ and let f_δ be the density function

$$f_\delta(u) := \delta \mathbf{1}_{[0, 0.5]}(u) + \mathbf{1}_{[0.5, 1]}(u)$$

For sufficiently small δ , a quantizer $x = (y_1, x_2)$ is a critical point of D^{2, f_δ} if and only if it satisfies

$$\begin{cases} 3x_1^2 - x_2^2 + 2y_1x_2 + 4y_1(\delta - 1) + 1 - \delta = 0 \\ x_1^2 - 3x_2^2 - 2y_1x_2 + 8x_2 - 4 = 0 \end{cases}$$

The preceding equation system has two solutions, $(0.25, 0.75)$ and $\left(\frac{5-9\delta}{8}, \frac{7-3\delta}{8}\right)$.

The quantizers $\left(\frac{5-9\delta}{8}, \frac{7-3\delta}{8}\right)$ is the absolute minimum of D^{2, f_δ} . The quantizer $(0.25, 0.75)$ is a stationary quantizer but one can easily verify that it is not a local minimum of D^{2, f_δ} .

0.12.2 Loss functions and log-concave density.

1. We show that the loss function used by Trushkin in [11] (section III) that is, even, convex, zero at 0, satisfy the assumptions (V) and (C) provided that μ satisfies (M). Let V be an even convex loss function satisfying (C) (i); the derivative of a convex function is right and left limited and bounded on compact sets. So, the function V satisfies (V). For all c_1, c_2 such that $-\infty < c_1 < c_2 < +\infty$, one has

$$|v(c-u)| \leq |v(c_1-u)| + |v(c_2-u)| \quad \forall u \in \mathbf{R}, \quad \forall c \in [c_1, c_2]$$

But, from lemma 1 in [11], one has for all $c \in (m, M)$,

$$\int v(c-u) f(u) du < +\infty$$

Thus

$$\sup_{c \in [c_1, c_2]} |v(c-u)| \in L^1(\mu)$$

Integrating by part then yields (see lemma 2, [11])

$$u \rightarrow V(c-u) f'(u) \in L^1(du) \quad \forall c \in (m, M)$$

which finally yields

$$\sup_{c \in [c_1, c_2]} |v(c-u)| f'(u) \in L^1(du)$$

As a result, (μ, V) satisfies the Assumption (C). So, Assumptions (V) and (C) contain Trushkin's convex losses.

2. Assumption (V) permit us to deal with non-convex loss function. For instance, let V be differentiable, even and Lipschitzian on compact sets so that V satisfies (V). Assume in addition that V satisfies (C) (i). Then (μ, V) satisfies Assumption (C).

Hence, Theorem 0.11 shows that the convexity of loss functions is not essential in the uniqueness problem.

Classical log-concave density functions are

1. constant and intervally supported density functions.
2. normal density functions $\mathcal{N}(m, \sigma^2)$

$$f(u) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(u-m)^2}{2\sigma^2}}$$

3. $\gamma(\alpha, \beta)$ density functions with $\alpha \leq 1$ and $\beta > 0$

$$f(u) = \frac{u^{\alpha-1} e^{-u/\beta}}{\beta^\alpha \Gamma(\alpha)} \mathbf{1}_{\mathbb{R}_+}(u)$$

In particular, exponential and $\chi^2(k)$, $k \geq 2$, density functions (the one degree χ^2 is log-convex).

4. $\beta(a, b)$ density functions with $a \geq 1$ and $b \geq 1$

$$f(u) = \frac{u^{a-1} (1-u)^{b-1}}{B_{a,b}} \mathbf{1}_{[0,1]}(u)$$

5. Every concave density functions.

Conclusion.

The main contribution of the geometrical approach to the uniqueness problem consists in two fact :

- First, we think that it provides a more intuitive setting than the Trushkin or Kieffer's approach since the Mountain Pass Lemma acts straight on local properties of the distortion.
- Second, it permit to extend the loss function used in previous works by showing in particular that the convexity of V is not essential in the uniqueness problem.

Moreover, it would provide an alternative to investigate new conditions on the density function f that ensure uniqueness. For instance, one may think about the log-convexity of the density function f and, in particular to the $\chi^2(1)$ density function.

Annex.

Proof of the Mountain Pass Lemma.

Here, we adapt the proof of the classical version of the Mountain Pass Lemma given in [10]. The classical proof deals essentially with local arguments. So, we can replace the coercivity of L by the property $(I_d - \varepsilon \nabla L)(O) \subset (O)$.

Proof. Assume there exists two strict minima for L , say y and z . Let $\beta := \inf_{p \in \mathcal{P}} \max_{x \in p} L(x)$ where \mathcal{P} is the set of compact and connected subsets of K containing y_1 and x_2 . Since K is connected, one has $\mathcal{P} \neq \emptyset$ and so, there is a minimising sequence (p_k) for β . Furthermore, one can define $p = \bigcup_{n \in \mathbb{N}} \bigcap_{k \geq n} p_k$. On one hand, the set p is compact and connected as decreasing limit of compact and connected sets. In addition, $\{y_1, x_2\} \subset p$ which implies $p \in \mathcal{P}$. As a result, $\max_{x \in p} L(x) \geq \beta$. On the other hand, for all n ,

$$\max_{x \in p} L(x) \leq \sup \left\{ L(x); x \in \bigcap_{k \geq n} p_k \right\}$$

but,

$$\sup \left\{ L(x); x \in \bigcap_{k \geq n} p_k \right\} = \sup_{k \geq n} \max_{p_k} L(x)$$

then,

$$\begin{aligned} \max_{x \in p} L(x) &\leq \sup_{k \geq n} \max_{p_k} L(x) \\ &= \beta \end{aligned}$$

Finally, $\max_{x \in p} L(x) = \beta$.

Let M be the compact set $\{x \in p; L(x) = \beta\}$; Assume that M does not contain any critical point; the continuity of ∇L on K yields the existence of $\delta > 0$ such that $\forall x \in M, |\nabla L(x)| > \delta$ and the existence of a K -open set U such that $|\nabla L(x)| > \frac{\delta}{2} \quad \forall x \in U$.

Let η be a continuous function $K \mapsto [0, 1]$ such that $\text{supp}(\eta) \subset U$ and such that $\eta \equiv 1$ on a K -open set containing M . Let $\phi : K \times \mathbf{R} \mapsto \mathbf{R}^N$ be the function

$$\phi(x, t) := x - t\eta(x) \nabla L(x).$$

The function L is continuously differentiable on O and ∇L admits a continuous extension on ∂K . So, by lemma 0.5, L is differentiable on K relatively to the topology induced on K , with gradient ∇L . But, for all $x \in K$, for all $t \in [0, \varepsilon]$, $\phi(x, t) \in K$. Indeed, $(I_d - \varepsilon \nabla L)(O) \subset O$ and the continuity of ∇L on K yields $(I_d - \varepsilon \nabla L)(K) \subset K$. So, by convexity of K , we obtain $\phi(K \times [0, \varepsilon]) \subset K$. Now, ϕ is differentiable on $[0, \varepsilon]$. So, by lemma 0.6, the function $t \mapsto L \circ \phi(x, t)$ is differentiable on $[0, \varepsilon]$ for all $x \in K$ and

$$\frac{d}{dt} L(\phi(x, t)) = -\eta(x) \langle \nabla L(\phi(x, t)), \nabla L(x) \rangle$$

Since for all $x \in U$, one has $|\nabla L(x)| \geq \frac{\delta}{2}$, we obtain

$$-\langle \nabla L(\phi(x, 0)), \nabla L(x) \rangle \leq \frac{\delta^2}{4} \quad \forall x \in U$$

By continuity of the function $(x, t) \mapsto \langle \nabla L(\phi(x, t)), \nabla L(x) \rangle$ on $K \times [0, \varepsilon]$ there exists $T \in]0, \varepsilon]$ such that

$$-\langle \nabla L(\phi(x, t)), \nabla L(x) \rangle \leq \frac{\delta^2}{8} \quad \forall (x, t) \in U \times [0, T]$$

Seeing that $\text{supp}(\eta) \subset U$, we obtain

$$\frac{d}{dt} L(\phi(x, t)) \leq -\eta(x) \frac{\delta^2}{8} \quad \forall (x, t) \in K \times [0, T]$$

Now, the function $t \mapsto L \circ \phi(x, t)$ is differentiable for every $x \in K$; in particular, we can say that

$$L(\phi(x, T)) = L(x) + \int_0^T \frac{d}{dt} L(\phi(x, t)) dt \quad (20)$$

$$\leq L(x) - \eta(x) \frac{T\delta^2}{2} \quad (21)$$

Let $p_T := \phi(p, T)$. The inequality (21) yields

$$\phi(x, T) \leq \beta - \frac{T}{2} \delta^2 \quad \text{for all } x \in M$$

$$\phi(x, T) \leq L(x) < \beta \quad \text{for all } x \in p \setminus M$$

So, one has

$$\max_{x \in p_T} L(x) < \beta$$

Finally, y_1 et x_2 are $\phi(\cdot, T)$ -invariant and since the function $L \circ \phi$ is continuous, the set p_T is compact and connected. As a consequence, $p_T \in \mathcal{P}$ which contradicts the minimaximality of the set p .

Then, one has $\{\nabla L = 0\} \cap M \neq \emptyset$. This set is K -closed and then also p -closed. Assume that all elements of p are K -relative minima. For all $y \in \{\nabla L = 0\} \cap M$, there exists a K -open set, say V , such that $\forall z \in V$, $L(z) \geq \beta$. So, one has $L(z) = \beta$ for all $z \in V \cap p$ which implies that all elements of $V \cap p$ are K -relative minima. Then, $V \cap p \subset \{\nabla L = 0\} \cap M$; so the set $\{\nabla L = 0\} \cap M$ is p -open. Since p is connected, we have $\{\nabla L = 0\} \cap M = p$, which contradicts the fact that $L(y_1) < \beta$. \diamond

Lemma 0.5 and 0.6.

We say that a function is K -differentiable if it is differentiable on K for the topology of the set K .

Lemma 0.5. *Let A, B be Banach spaces. Let $K \subset A$ be a convex set with non-empty interior. Let $L : K \rightarrow B$ be a continuous function on K , continuously differentiable on $\overset{\circ}{K}$. Assume that DL admits a continuous extension on ∂K .*

Then, L is continuously K -differentiable on K and DL is uniquely determined on ∂K .

Proof. We have to show that for all $x \in \partial K$, $\frac{\|L(x+h) - L(x) - DL(x)h\|}{\|h\|} \rightarrow 0$ as $h \rightarrow 0$ with $h \in K - x$. Let $x \in \partial K$ and $\varepsilon > 0$. Since K is a convex set, there exist $y \in \overset{\circ}{K} \cap B(x, \|h\|)$ such that $\|x+h-y\| \leq \|h\|^2$. Moreover, one has $]x, y[\subset \overset{\circ}{K}$, $]x+h, y[\subset \overset{\circ}{K}$ and $\|x-y\| \leq \|h\|$.

But,

$$\begin{aligned} \|L(x+h) - L(x) - DL(x)h\| &\leq \|L(x+h) - L(y)\| + \|L(y) - L(x) - DL(x)(y-x)\| \\ &\quad + \|DL(x)h - DL(x)(y-x)\| \\ &\leq 2 \sup_{t \in]x+y,[} \|DL(t)\| \|h\|^2 + \sup_{t \in]x,y,[} \|DL(t) - DL(x)\| \|h\| \\ &\quad + 2 \|DL(x)\| \|h\|^2 \end{aligned}$$

The function DL is continuous on K ; so, for small enough $\|h\|$, we obtain

$$\|L(x+h) - L(x) - DL(x)h\| \leq \varepsilon \|h\|$$

The differential of L is uniquely determined on ∂K for $(K-x) \cap B(0, \varepsilon)$ contains a basis of A for all $x \in \partial K$ and for all $\varepsilon > 0$ \diamond

Lemma 0.6. *Let A, B, C be Banach spaces. Let $K \subset A$ be a convex set with non-empty interior. Let $\phi : A \rightarrow K$ be differentiable and $L : K \rightarrow C$ be K -differentiable on K . Then, $L \circ \phi$ is differentiable on A .*

Bibliography

- [1] C. Bouton, G. Pagès, "Self-organization and a.s. Convergence of the One-dimensional Kohonen Algorithm with Non-uniformly Distributed Stimuli", *Stochastic Processes and their Applications* 47(1993) 249-274, North-Holland.
- [2] P.E. Fleischer, "Sufficient Conditions for Achieving Minimum Distortion in a Quantizer", *IEEE Int. Conv. Rec.*, 1964, part 1, pp 104-111.
- [3] J.C. Fort, G. Pagès, "About the a.s. Convergence of the Kohonen Algorithm with a General Neighborhood Function", *The annals of applied probability*, vol 5, n°4, 1995.
- [4] R.M Gray, E.D. Karnin, "Multiple Local Optima in Vector Quantizers", *IEEE Transactions on Information Theory*, vol IT28, n°2, march 1982, pp 256-261.
- [5] J.C. Kieffer, "Uniqueness of Locally Optimal Quantizer for Log-concave Density and Convex Error Weighting Function", *IEEE Transactions on Information Theory*, vol IT29, n°1, january 1983, pp 42-47.
- [6] J.C. Kieffer, "Exponential Rate of Convergence for Lloyd's Method I", *IEEE Transactions on Information Theory*, vol IT28, n°2, march 1982, pp 205-210.
- [7] D. Lambertson, G. Pagès, "On the Critical Points of the 1-dimensional Competitive Learning Vector Quantization Algorithm", Proceedings of the **ESANN'96**, M. Verleysen d., D Facto editeur, Bruxelles, pp97-102, 1996.
- [8] S.P. Lloyd, "Least Squares Quantization in P.C.M.", *IEEE Transactions on Information Theory*, vol IT28, n°2, march 1982, pp 129-137.
- [9] M. Struwe, "Variational Methods (Applications to Non Linear p.d. & Hamiltonian Systems)", Springer, chapitre 2.
- [10] A.V. Trushkin, "Sufficient Conditions for Uniqueness of a Locally Optimal Quantizer for a Class of Convex Error Weighting Functions", *IEEE Transactions on Information Theory*, vol IT28, n°2, march 1982, pp 187-198.
- [11] A.V. Trushkin, "Monotony of Lloyd's Method II for Log-concave Density and Convex Error Weighting Function", *IEEE Transactions on Information Theory*, vol IT30, n°2, march 1984, pp 380-383.
- [12] A.V. Trushkin, "On the Design of an Optimal Quantizer", *IEEE Transactions on Information Theory*, vol 39, n°4, july 1993, pp 1180-1194.
- [13] R. Bhatia, **Matrix Analysis**, Graduate Texts in Mathematics 169. Springer.