



**HAL**  
open science

# Intrinsic stationarity for vector quantization: Foundation of dual quantization

Gilles Pagès, Benedikt Wilbertz

► **To cite this version:**

Gilles Pagès, Benedikt Wilbertz. Intrinsic stationarity for vector quantization: Foundation of dual quantization. 2010. hal-00528485v1

**HAL Id: hal-00528485**

**<https://hal.science/hal-00528485v1>**

Preprint submitted on 21 Oct 2010 (v1), last revised 26 Mar 2012 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Intrinsic stationarity for vector quantization: Foundation of dual quantization

Gilles Pagès      Benedikt Wilbertz

October 22, 2010

## Abstract

We develop a new approach to vector quantization, which guarantees an intrinsic stationarity property that also holds, in contrast to regular quantization, for non-optimal quantization grids. This goal is achieved by replacing the usual nearest neighbor projection operator for Voronoi quantization by a random splitting operator, which maps the random source to the vertices of a triangle of  $d$ -simplex. In the quadratic Euclidean case, it is shown that these triangles or  $d$ -simplices make up a Delaunay triangulation of the underlying grid.

Furthermore, we prove the existence of an optimal grid for this Delaunay – or dual – quantization procedure. We also provide a stochastic optimization method to compute such optimal grids, here for higher dimensional uniform and normal distributions. A crucial feature of this new approach is the fact that it automatically leads to a second order quadrature formula for computing expectations, regardless of the optimality of the underlying grid.

*Keywords:* Quantization, Stationarity, Voronoi tessellation, Delaunay triangulation, Numerical integration.

*MSC:* 60F25, 65C50, 65D32

## 1 Introduction and motivation

Quantization of random variables aims at finding the best  $p$ -th mean approximation to a r.v.  $X : (\Omega, \mathcal{S}, \mathbb{P}) \rightarrow (\mathbb{R}^d, \mathcal{B}^d)$  and  $\mathbb{R}^d$  equipped with some norm  $\|\cdot\|$ . That means, for  $X \in L^p(\mathbb{P})$ ,  $p \geq 1$  that we have to minimize

$$\mathbb{E} \min_{x \in \Gamma} \|X - x\|^p \tag{1}$$

over all grids  $\Gamma \subset \mathbb{R}^d$  of a given size. This problem has its origin in the fields of signal processing in the late 1940s. A mathematically rigorous and comprehensive exposition of this topic can be found in the book [7].

Using the nearest neighbor projection, we are able to construct a random variable  $\widehat{X}^\Gamma$ , which achieves the minimum in (1). Such an approximation, which is called Voronoi quantization, has been successfully applied to various problems in applied probability theory and mathematical finance, *e.g.* multi-asset American/Bermudan style options pricing and  $\delta$ -hedging (see [1, 2]), swing options, supply gas contract, on energy markets (Stochastic control) (see [3, 4, 5]), non-linear filtering method for stochastic volatility estimation (see [9, 12, 14, 15]), discretization of SPDE's (stochastic Zakai and McKean-Vlasov equations) (see [6]).

Especially we may use optimal quantizations to establish numerical cubature formulas, *i.e.* to approximate  $\mathbb{E}F(X)$  by

$$\mathbb{E}F(\widehat{X}^\Gamma) = \sum_{x \in \Gamma} w_x \cdot F(x),$$

where  $w_x = \mathbb{P}(\widehat{X}^\Gamma = x)$ .

Such a cubature formula is known to be optimal in the class of Lipschitz functionals and it holds for a Lipschitz functional  $F$

$$|\mathbb{E}F(X) - \mathbb{E}F(\widehat{X}^\Gamma)| \leq [F]_{\text{Lip}} \mathbb{E}\|X - \widehat{X}^\Gamma\|. \quad (2)$$

If  $F$  exhibits a bit more smoothness, *i.e.* is piecewise differentiable with Lipschitz derivative and  $\widehat{X}$  fulfills the so-called *stationarity property*

$$\mathbb{E}(X|\widehat{X}^\Gamma) = \widehat{X}^\Gamma, \quad (3)$$

we can derive by means of a Taylor expansion the second order rate

$$|\mathbb{E}F(X) - \mathbb{E}F(\widehat{X}^\Gamma)| \leq [F']_{\text{Lip}} \mathbb{E}\|X - \widehat{X}^\Gamma\|^2.$$

Unfortunately, the stationarity property for the Voronoi quantization  $\widehat{X}^\Gamma$  is a rather fragile object, since it only holds for grids  $\Gamma$  which are especially tailored and optimized for the distribution of  $X$ .

That means, that if a grid  $\Gamma$ , which has been originally constructed and optimized for  $X$ , is employed to approximate a r.v.  $Y$  which only slightly differs from  $X$ , then  $\Gamma$  might be still an arbitrary good quantization for  $Y$ , *i.e.*  $\mathbb{E}\|Y - \widehat{Y}^\Gamma\|^p$  is very close to the optimal quantization error, but the stationarity property (3) is in general violated. Thus, only the first order bound (2) is in this case valid for a cubature formula based on a Voronoi quantization of  $Y$ .

In this paper, we look for an alternative to the nearest neighbor projection operator and the Voronoi quantization, which will be capable of preserving some stationarity property in the above setting. In order to achieve this, we pass on to a product space  $(\Omega_0 \times \Omega, \mathcal{S}_0 \otimes \mathcal{S}, \mathbb{P}_0 \otimes \mathbb{P})$  and introduce a *random splitting operator*  $\mathcal{J}_\Gamma : \Omega_0 \times \mathbb{R}^d \rightarrow \Gamma$ , which satisfies

$$\mathbb{E}(\mathcal{J}_\Gamma(Y)|Y) = Y$$

for any  $\mathbb{R}^d$ -valued r.v.  $Y$  defined on  $(\Omega, \mathcal{S}, \mathbb{P})$  such that  $\text{supp}(\mathbb{P}_Y) \subset \text{conv}(\Gamma)$ . As a matter of facts, such an operator fulfills the so-called *intrinsic stationarity property*

$$\mathbb{E}(\mathcal{J}_\Gamma(\xi)) = \xi, \quad \xi \in \text{conv}(\Gamma). \quad (4)$$

Although this stationarity differs from the one defined above, one may again derive a second order error bound for a differentiable function  $F$  with Lipschitz derivative

$$|\mathbb{E}F(Y) - \mathbb{E}F(\mathcal{J}_\Gamma(Y))| \leq [F']_{\text{Lip}} \mathbb{E}\|Y - \mathcal{J}_\Gamma(Y)\|^2,$$

which now holds for any r.v.  $Y$  regardless of the grid  $\Gamma$ , except satisfying  $\text{supp}(\mathbb{P}_Y) \subset \text{conv}(\Gamma)$ .

One may naturally ask at this stage for the best possible approximation power of  $\mathcal{J}_\Gamma(X)$  to  $X$ , *i.e.* minimize the  $p$ -th power mean error

$$\mathbb{E}\|X - \mathcal{J}_\Gamma(X)\|^p$$

over all grids of size not exceeding  $n$  and all random operators  $\mathcal{J}_\Gamma$ , which fulfill the intrinsic stationarity property (4).

This means, that we will deal for  $n \in \mathbb{N}$  with the mean error modulus

$$d_n^p(X) = \inf \left\{ \mathbb{E} \|X - \mathcal{J}_\Gamma(X)\|^p : \Gamma \subset \mathbb{R}^d, |\Gamma| \leq n, \text{supp}(\mathbb{P}_X) \subset \text{conv}(\Gamma), \right. \\ \left. \mathcal{J}_\Gamma : \Omega_0 \times \mathbb{R}^d \rightarrow \Gamma \text{ intrinsic stationary} \right\}. \quad (5)$$

It will turn out in Section 2 that the problem of finding an optimal random operator  $\mathcal{J}_\Gamma$  for a grid  $\Gamma = \{x_1, \dots, x_k\}, k \leq n$ , is equivalent to solving the Linear Programming problem

$$\min_{\lambda \in \mathbb{R}^k} \sum_{i=1}^k \lambda_i \|X(\omega) - x_i\|^p. \quad (6) \\ \text{s.t. } \begin{bmatrix} x_1 & \dots & x_k \\ 1 & \dots & 1 \end{bmatrix} \lambda = \begin{bmatrix} X(\omega) \\ 1 \end{bmatrix}, \lambda \geq 0$$

Defining the local dual quantization function as

$$F^p(\xi, \Gamma) = \min_{\lambda \in \mathbb{R}^k} \sum_{i=1}^k \lambda_i \|\xi - x_i\|^p, \\ \text{s.t. } \begin{bmatrix} x_1 & \dots & x_k \\ 1 & \dots & 1 \end{bmatrix} \lambda = \begin{bmatrix} \xi \\ 1 \end{bmatrix}, \lambda \geq 0$$

we will show that

$$d_n^p(X) = \inf \{ \mathbb{E} F^p(X; \Gamma) : \Gamma \subset \mathbb{R}^d, |\Gamma| \leq n \}. \quad (7)$$

This means, that the dual quantization problem actually consists of two phases: during the first one we have to locally solve the optimization problem (6), whereas phase two, which consists of the global optimization over all possible grids in (7), is the more involved problem. It is highly non-linear and contains a probabilistic component by contrast to phase one which can be considered more or less as deterministic.

Moreover, we will see in section 3 that the solution to the Linear Programming (6) is in the quadratic Euclidean case completely determined by the Delaunay triangulation spanned by  $\Gamma$  and this structure is, in the graph theoretic sense, the dual counterpart of the Voronoi diagram, on which regular quantization is based. That is actually also the reason, why we call this new approach dual or Delaunay quantization.

In section 2, we moreover give a generalization of the dual quantization idea to non-compactly supported random variables. For those and the compactly supported r.v.'s we prove the existence of optimal quantizers in section 4, *i.e.* the fact, that there are sets  $\Gamma$ , which actually achieve the infimum in (5). Finally, in section 5, we give numerical illustrations of some optimal dual quantizers and numerical procedures to generate them.

In a companion paper [11], we establish the counterpart of the celebrated Zador theorem for regular vector quantization: namely we elucidate the sharp rate for the mean dual quantization error modulus defined in section 2 below.

We also provide in [11] a non-asymptotic version of this theorem, which corresponds to Pierce's Lemma.

NOTATION:  $\bullet u^T$  will denote the transpose of the column vector  $u \in \mathbb{R}^d$ .

$\bullet$  Let  $u = (u_1, \dots, u_d) \in \mathbb{R}^d$ , we write  $u \geq 0$  (resp.  $> 0$ ) if  $u_i \geq 0$  (resp.  $> 0$ )  $\forall i = 1, \dots, d$ .

## 2 Dual quantization and intrinsic stationarity

First, we briefly recall the definition of the ‘‘regular’’ vector quantization problem for a r.v.  $X : (\Omega, \mathcal{S}, \mathbb{P}) \rightarrow (\mathbb{R}^d, \mathcal{B}^d)$  and  $\mathbb{R}^d$  equipped with a norm  $\|\cdot\|$ .

**Definition 1.** Let  $X \in L_{\mathbb{R}^d}^p(\mathbb{P})$  for some  $p \in [1, +\infty)$ .

1. We define the (regular)  $L^p$ -mean quantization error for a grid  $\Gamma = \{x_1, \dots, x_k\} \subset \mathbb{R}^d$  as

$$e_p(X; \Gamma) = \left\| \min_{1 \leq i \leq k} \|X - x_i\| \right\|_{L^p} = \left( \mathbb{E} \min_{1 \leq i \leq k} \|X - x_i\|^p \right)^{1/p},$$

2. The optimal regular quantization error, which can be achieved by a grid  $\Gamma$  of size not exceeding  $n \in \mathbb{N}$ , is given by

$$e_{n,p}(X) = \inf \{ e_p(X; \Gamma) : \Gamma \subset \mathbb{R}^d, |\Gamma| \leq n \}.$$

*Remark.* Since we will frequently consider the  $p$ -th power of  $e_p(X; \Gamma)$  and  $e_{n,p}(X)$ , we will drop a duplicate index  $p$  and write, e.g.  $e_n^p(X)$  instead of  $e_{n,p}^p(X)$ .

It can be shown, that (at least) one optimal quantizer actually exists, i.e. for every  $n \in \mathbb{N}$  there is a grid  $\Gamma \subset \mathbb{R}^d$  with  $|\Gamma| \leq n$  such that

$$e_p(X; \Gamma) = e_{n,p}(X).$$

Moreover, this definition of the optimal quantization error is in fact equivalent to defining  $e_n^p(X)$  as the best approximation error which can be achieved by a Borel transformation or by a discrete r.v.  $\widehat{X}$  taking at most  $n$  values:

**Proposition 1.** *Let  $X \in L^p(\mathbb{P})$ ,  $n \in \mathbb{N}$ . Then*

$$\begin{aligned} e_n^p(X) &= \inf \{ \mathbb{E} \|X - f(X)\|^p : f : \mathbb{R}^d \rightarrow \mathbb{R} \text{ Borel mb, } |f(\mathbb{R}^d)| \leq n \} \\ &= \inf \{ \mathbb{E} \|X - \widehat{X}\|^p : \widehat{X} \text{ is a r.v. with } |\widehat{X}(\Omega)| \leq n \}. \end{aligned}$$

The proof of this Proposition is based on the construction of a Voronoi quantization of a r.v. by means of the nearest neighbor projection.

Therefore, let  $\Gamma = \{x_1, \dots, x_k\} \subset \mathbb{R}^d$  be a grid and denote by  $(C_i(\Gamma))_{1 \leq i \leq k}$  a Borel partition of  $\mathbb{R}^d$  satisfying

$$C_i(\Gamma) \subset \{ \xi \in \mathbb{R}^d : \|\xi - x_i\| \leq \min_{1 \leq j \leq k} \|\xi - x_j\| \}.$$

Such a partition is called a *Voronoi partition* generated by  $\Gamma$  and we may define the corresponding *nearest neighbor projection* as

$$\pi_\Gamma(\xi) = \sum_{1 \leq i \leq k} x_i \mathbb{1}_{C_i(\Gamma)}(\xi).$$

The discrete r.v.

$$\widehat{X}^{\Gamma, \text{Vor}} = \pi_\Gamma(X) = \sum_{1 \leq i \leq k} x_i \mathbb{1}_{C_i(\Gamma)}(X).$$

is called *Voronoi Quantization* induced by  $\Gamma$  and satisfies

$$e^p(X; \Gamma) = \mathbb{E} \|X - \pi_\Gamma(X)\|^p.$$

As already mentioned in the introduction, the concept of stationarity plays an important role in the application of quantization. A quantization  $\widehat{X}$  is said to be *stationary* for the r.v.  $X$ , if it satisfies

$$\mathbb{E}(X|\widehat{X}) = \widehat{X}. \quad (8)$$

It is well known that in the quadratic Euclidean case, i.e.  $p = 2$  and  $\|\cdot\|$  is the Euclidean norm, any optimal quantization, that is a r.v.  $\widehat{X}$  with  $|\widehat{X}(\Omega)| \leq n$  for an  $n \in \mathbb{N}$  and  $\mathbb{E} \|X - \widehat{X}\|^p = e_n^p(X)$ , fulfills this property.

Moreover, this stationarity condition is equivalent to the first order optimality criterion of the optimization problem

$$\mathbb{E} \min_{1 \leq i \leq n} |X - x_i|^2 \rightarrow \min_{x_1, \dots, x_n \in \mathbb{R}^d},$$

i.e. the Voronoi quantization  $\widehat{X}^{\Gamma, \text{Vor}}$  of a grid  $\Gamma = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$  satisfies the stationarity property (8) for a r.v.  $X$ , whenever  $\Gamma$  is a zero of the first order derivative of the mapping  $(x_1, \dots, x_n) \mapsto \mathbb{E} \min_{1 \leq i \leq n} |X - x_i|^2$ .

By means of this stationarity property (8), we can derive the following second order error bound for a cubature formula based on quantization.

**Proposition 2.** *Let  $X \in L^2(\mathbb{P})$  and assume that  $F \in C^{1,1}(\mathbb{R}^d)$  is differentiable with Lipschitz derivative. If the quantization  $\widehat{X}^\Gamma$  for a grid  $\Gamma = \{x_1, \dots, x_n\} = \widehat{X}^\Gamma(\Omega)$ ,  $n \in \mathbb{N}$  satisfies*

$$\mathbb{E}(X|\widehat{X}^\Gamma) = \widehat{X}^\Gamma,$$

then it holds for the cubature formula  $\mathbb{E}F(\widehat{X}^\Gamma) = \sum_{i=1}^n \mathbb{P}(\widehat{X}^\Gamma = x_i) \cdot F(x_i)$

$$|\mathbb{E}F(X) - \mathbb{E}F(\widehat{X}^\Gamma)| \leq [F']_{\text{Lip}} \mathbb{E}\|X - \widehat{X}^\Gamma\|^2.$$

*Proof.* From a Taylor expansion we obtain for  $\widehat{X} = \widehat{X}^\Gamma$

$$|F(X) - F(\widehat{X}) - F'(\widehat{X})(X - \widehat{X})| \leq [F']_{\text{Lip}} \|X - \widehat{X}\|^2,$$

so that taking conditional expectations and applying Jensen's inequality yield

$$|\mathbb{E}(F(X)|\widehat{X}) - F(\widehat{X}) - \mathbb{E}(F'(\widehat{X})(X - \widehat{X})|\widehat{X})| \leq [F']_{\text{Lip}} \mathbb{E}(\|X - \widehat{X}\|^2|\widehat{X}).$$

The stationarity assumption then implies

$$\mathbb{E}(F'(\widehat{X})(X - \widehat{X})|\widehat{X}) = F'(\widehat{X}) \mathbb{E}((X - \widehat{X})|\widehat{X}) = 0,$$

so that the assertion follows again from taking expectations and Jensen's inequality.  $\square$

Unfortunately, the above stationarity is a rather fragile property, since it only holds for Voronoi quantizations, whose underlying grid is specifically optimized for the distribution of  $X$ . Thus, this stationarity will in general fail, as soon as we modify the underlying r.v. even only slightly. Nevertheless, there is a second way to derive the second order error bound of Proposition 2: Assume that  $\widehat{X}$  is a discrete r.v. satisfying a somewhat dual stationarity property

$$\mathbb{E}(\widehat{X}|X) = X. \tag{9}$$

In this case we can perform, as in the proof of Proposition 2, a Taylor expansion, but this time with respect to  $X$ . We then conclude from (9)

$$\mathbb{E}(F'(X)(X - \widehat{X})|X) = 0$$

so that finally the same assertion will hold.

As we will see later on, this stationarity condition will be intrinsically fulfilled by the dual quantization operator. Thus, this new approach will be very robust with respect to changes in the underlying r.v.s, since it always preserves stationarity.

## 2.1 Definition of dual quantization

We define here the dual quantization error by means of the local dual quantization error  $F_p$ , since, doing so, we are able to introduce dual quantization along the lines of regular quantization. The stationarity property (9) will then appear as characterizing property of the Delaunay quantization and the dual quantization operator, the counterpart of Voronoi quantization and the nearest neighbor projection.

The equivalence of the Definition 2 and (5) will be given in Theorem 2, which provides an analog statement for dual quantization to Proposition 1.

Without loss of generality assume from here on that

$$\text{span}(\text{supp}(\mathbb{P}_X)) = \mathbb{R}^d,$$

i.e.  $X$  is a true  $d$ -dimensional random variable. Otherwise we would reduce  $d$ .

**Definition 2.** Let  $X \in L^p(\mathbb{P})$  for some  $p \in [1, \infty)$ .

1. We define the local dual quantization error for a grid  $\Gamma = \{x_1, \dots, x_k\} \subset \mathbb{R}^d$  as

$$F_p(\xi; \Gamma) = \inf \left\{ \left( \sum_{1 \leq i \leq k} \lambda_i \|\xi - x_i\|^p \right)^{1/p} : \lambda_i \in [0, 1] \text{ and } \sum_{1 \leq i \leq k} \lambda_i x_i = \xi, \sum_{1 \leq i \leq k} \lambda_i = 1 \right\}$$

2. The  $L^p$ -mean dual quantization error for  $X$  induced by the grid  $\Gamma$  is then given by

$$d_p(X; \Gamma) = \|F_p(X; \Gamma)\|_{L^p} = \left( \mathbb{E} \inf \left\{ \sum_{1 \leq i \leq k} \lambda_i \|X - x_i\|^p : \lambda_i \in [0, 1], \sum_{1 \leq i \leq k} \lambda_i x_i = X, \sum_{1 \leq i \leq k} \lambda_i = 1 \right\} \right)^{1/p}$$

3. The optimal dual quantization error, which can be achieved by a grid  $\Gamma$  of size not exceeding  $n$  will be denoted by

$$d_{n,p}(X) = \inf \{ d_p(X; \Gamma) : \Gamma \subset \mathbb{R}^d, |\Gamma| \leq n \}.$$

*Remark.* 1. Note that, as in the case of regular quantization, the optimal dual quantization error depends actually only on the distribution of  $X$ .

2. In many cases we will deal with the  $p$ -th power of  $F_p$ ,  $d_p$  and  $d_{n,p}$ . To avoid duplicating indices, we will write  $F^p$ ,  $d^p$  and  $d_{n,p}^p$  instead of  $F_p^p$ ,  $d_p^p$  and  $d_{n,p}^p$ .

Denoting  $\Gamma = \{x_1, \dots, x_k\}$ , we recognize that  $F^p(\xi; \Gamma)$  is given by the linear programming problem

$$\begin{aligned} \min_{\lambda \in \mathbb{R}^k} \sum_{i=1}^k \lambda_i \|\xi - x_i\|^p & \quad (\text{LP}) \\ \text{s.t. } \begin{bmatrix} x_1 & \dots & x_k \\ 1 & \dots & 1 \end{bmatrix} \lambda &= \begin{bmatrix} \xi \\ 1 \end{bmatrix}, \lambda \geq 0 \end{aligned}$$

Clearly, we have  $F^p(\xi; \Gamma) \geq 0 \forall \xi \in \mathbb{R}^d, \Gamma \subset \mathbb{R}^d$ , so that it follows from the constraints

$$\begin{bmatrix} x_1 & \dots & x_k \\ 1 & \dots & 1 \end{bmatrix} \lambda = \begin{bmatrix} \xi \\ 1 \end{bmatrix}, \quad \lambda \geq 0 \quad (10)$$

that (LP) has a finite solution if and only if  $\xi \in \text{conv}\{\Gamma\}$ .

**Proposition 3.** (a) Let  $p \in [1, +\infty)$  and assume  $\text{supp}(\mathbb{P}_X)$  is compact. For every  $n \geq d + 1$ ,  $d_{n,p}(X) < +\infty$  (and for every  $n \in \{1, \dots, d\}$ ,  $d_{n,p}(X) = +\infty$ ).

(b) Let  $p \in (1, +\infty)$ . It holds

$$\{d_{n,p}(X; \cdot) < +\infty\} = \{\Gamma \subset \mathbb{R}^d : \text{conv}(\Gamma) \supset \text{supp}(\mathbb{P}_X)\}.$$

*Proof.* (a) Let  $\xi_0 \in \text{supp}(\mathbb{P}_X)$  and  $R > 0$  such that  $\text{supp} \mathbb{P}_X \subset B_{\ell^\infty}(\xi_0, \frac{R}{2})$  (closed ball w.r.t. the  $\ell^\infty$ -norm). Since  $[-\frac{R}{2}, \frac{R}{2}]^d \subset -\frac{R}{2}\mathbf{1} + R\mathcal{S}_d$  where  $\mathcal{S}_d$  denotes the canonical simplex. Consequently

$$\text{supp}(\mathbb{P}_X) \subset \xi_0 - \frac{R}{2}\mathbf{1} + R\mathcal{S}_d = \text{conv}(\Gamma_0), \quad \Gamma_0 = \{\xi_0 - R/2 + Re^j, j = 0, \dots, d\}$$

where  $e^0 = 0$  and  $(e^j)_{1 \leq j \leq d}$  denotes the canonical basis of  $\mathbb{R}^d$ . Consequently

$$\forall \xi \in \text{supp}(\mathbb{P}_X), \quad F_p(\xi; \Gamma_0) \leq \delta(\Gamma_0)$$

where  $\delta(A) := \sup_{x, y \in A} \|x - y\|$ . More generally, for every grid  $\Gamma$  such that  $\text{supp}(\mathbb{P}_X) \subset \text{conv}(\Gamma)$ ,  $F_p(\xi; \Gamma) < +\infty$  for every  $\xi \in \text{supp} \mathbb{P}_X$ .

Hence, for every  $n \geq |\Gamma_0| = d + 1$ ,

$$d_{n,p}(X) \leq \delta(\Gamma_0).$$

If  $n \leq d$ , the convex hull of a grid  $\Gamma$  cannot contain  $\text{supp}(\mathbb{P}_X)$ : if so it contains its convex hull as well which is impossible since it has a nonempty interior whereas the dimension of  $\text{conv}(\Gamma)$  is at most  $n - 1$ -dimensional.

(b) It follows from what precedes that  $d_{n,p}(X; \Gamma) < +\infty$  if  $\text{conv}(\Gamma) \supset \text{supp}(\mathbb{P}_X)$ . Conversely, if  $\text{conv}(\Gamma) \not\supset \text{supp}(\mathbb{P}_X)$ , there exists  $\xi_0 \in \text{supp}(\mathbb{P}_X) \setminus \text{conv}(\Gamma)$ . Let  $\varepsilon_0 > 0$  such that  $B(\xi_0, \varepsilon_0) \cap \text{conv}(\Gamma) = \emptyset$ . On  $B(\xi_0, \varepsilon_0)$ ,  $F_p(\cdot, \Gamma) \equiv +\infty$  and  $\mathbb{P}_X(B(\xi_0, \varepsilon_0)) > 0$ , hence  $d_{n,p}(X; \Gamma) = +\infty$ .  $\square$

## 2.2 Preliminaries on the local dual quantization functional

Before we deal in detail with the dual quantization error for random variables, we have to derive some basic properties for the local dual quantization error functional  $F_p$ .

To ease notations, we may introduce whenever  $\Gamma$  and/or  $\xi$  are fixed the abbreviations

$$A = \begin{bmatrix} x_1 & \cdots & x_k \\ 1 & \cdots & 1 \end{bmatrix}, \quad b = \begin{bmatrix} \xi \\ 1 \end{bmatrix}, \quad c = \begin{bmatrix} \|\xi - x_1\|^p \\ \vdots \\ \|\xi - x_k\|^p \end{bmatrix}$$

so that the (LP) can be written as

$$\begin{aligned} & \min_{\lambda \in \mathbb{R}^k} \lambda^T c. \\ & \text{s.t. } A\lambda = b, \lambda \geq 0 \end{aligned}$$

Moreover, for any set  $I \subset \{1, \dots, k\}$  we denote by  $A_I = [a_{ij}]_{j \in I}$  the submatrix of  $A$  with columns corresponding to the indices in  $I$ , and by  $c_I = [c_i]_{i \in I}$  will be the subvector of  $c$  which rows are determined by  $I$ .

Since it follows from Proposition 3 that any grid  $\Gamma \subset \mathbb{R}^d$  with  $\text{aff. dim}\{\Gamma\} < d$  yields  $d_p(X; \Gamma) = +\infty$ , we will restrict in the sequel to grids with  $\text{aff. dim}\{\Gamma\} = d$ , which is equivalent to  $\text{rk} \begin{pmatrix} x_1 & \cdots & x_k \\ 1 & \cdots & 1 \end{pmatrix} = d + 1$ .

It is then classical background, that, for every  $\xi \in \text{conv}(\Gamma)$ , (LP) has a solution  $\lambda^* \in \mathbb{R}^k$ , which is given by an extremal point of the compact set of linear constraints (10). In terms of Linear Programming theory, this corresponds to the existence of a fundamental basis  $I^* \subset \{1, \dots, k\}$ , such that  $|I^*| = d + 1$ , the columns  $[x_j]_1, j \in I^*$  are linearly independent and after reordering the rows we have

$$\lambda^* = \begin{bmatrix} A_{I^*}^{-1} b \\ 0 \end{bmatrix},$$

which means, that the columns of  $\lambda^*$  corresponding to  $I$  are given by  $A_{I^*}^{-1} b$ , the remaining ones are equal to 0.



Consequently, the Linear Program (LP) always admits a solution  $\lambda^*$ , whose non-zero components correspond to at most  $d + 1$  affinely independent points  $x_j$  in  $\Gamma$ , *i.e.* an optimal triangle in  $\mathbb{R}^2$  or  $d$ -simplex in  $\mathbb{R}^d$ .

Since the whole minimization problem can therefore be restricted to such triangles or  $d$ -simplices, we introduce the set of bases (or admissible indices) for a grid  $\Gamma = \{x_1, \dots, x_k\} \subset \mathbb{R}^d$  as

$$\mathcal{I}(\Gamma) = \{I \subset \{1, \dots, k\} : |I| = d + 1 \text{ and } \text{rk}(A_I) = d + 1\}.$$

Moreover, we denote the optimality region for a basis  $I \in \mathcal{I}(\Gamma)$  by

$$D_I(\Gamma) = \left\{ \xi \in \mathbb{R}^d : \lambda_I^* = A_I^{-1} \begin{bmatrix} \xi \\ 1 \end{bmatrix} \geq 0 \text{ and } \sum_{j \in I} \lambda_j^* \|\xi - x_j\|^p = \min_{\lambda \in \mathbb{R}^k} \sum_{i=1}^k \lambda_i \|\xi - x_i\|^p \right\}.$$

s.t.  $\begin{bmatrix} x_1 & \dots & x_k \\ 1 & \dots & 1 \end{bmatrix} \lambda = \begin{bmatrix} \xi \\ 1 \end{bmatrix}, \lambda \geq 0$

A useful reformulation of the above linear programming problem is given by its dual version:

**Proposition 4** (Duality). *The dual problem of (LP) reads*

$$\begin{aligned} \min_{\lambda \in \mathbb{R}^k} \sum_{i=1}^k \lambda_i \|\xi - x_i\|^p &= \max_{(u_1, u_2) \in \mathbb{R}^{d+1}} u_1^T \xi + u_2 \\ \text{s.t. } \begin{bmatrix} x_1 & \dots & x_k \\ 1 & \dots & 1 \end{bmatrix} \lambda &= \begin{bmatrix} \xi \\ 1 \end{bmatrix}, \lambda \geq 0 \end{aligned} \quad \text{s.t. } \begin{bmatrix} x_1^T & 1 \\ \vdots & \vdots \\ x_k^T & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \leq \begin{bmatrix} \|\xi - x_1\|^p \\ \vdots \\ \|\xi - x_k\|^p \end{bmatrix} \quad \text{(DLP)}$$

$$= \max_{u \in \mathbb{R}^d} \min_{1 \leq i \leq k} \{ \|\xi - x_i\|^p + u^T (\xi - x_i) \}.$$

*Proof.* Follows from classical duality for Linear Programs (see e.g. [8]). □

An important criterion to check, whether a triangle or a  $d$ -simplex in  $\Gamma$  is optimal, is given by the following characterization of optimality in Linear Programs:

**Proposition 5** (Optimality Conditions). (a) *If a basis  $I \in \mathcal{I}(\Gamma)$  is primal feasible, *i.e.**

$$\lambda_I = A_I^{-1} \begin{bmatrix} \xi \\ 1 \end{bmatrix} \geq 0,$$

*as well as dual feasible, *i.e.**

$$A^T u \leq \begin{bmatrix} \|\xi - x_1\|^p \\ \vdots \\ \|\xi - x_k\|^p \end{bmatrix} \quad \text{for } u = (A_I^T)^{-1} c_I,$$

*then*

$$\sum_{j \in I} \lambda_j \|\xi - x_j\|^p = \begin{bmatrix} \xi \\ 1 \end{bmatrix}^T u,$$

$\lambda$  and  $u$  are optimal for (LP) resp. (DLP) and  $I$  is called optimal basis.

(b) *Conversely, if  $I \in \mathcal{I}(\Gamma)$  is an optimal basis, which is additionally non-degenerated for (LP), *i.e.* there exist  $\lambda \in \mathbb{R}^k$  and  $u \in \mathbb{R}^{d+1}$  such that  $\lambda_I = A_I^{-1} \begin{bmatrix} \xi \\ 1 \end{bmatrix} > 0$  and*

$$\sum_{j \in I} \lambda_j \|\xi - x_j\|^p = \begin{bmatrix} \xi \\ 1 \end{bmatrix}^T u,$$

*then it holds*

$$A_I^T u = c_I.$$

*Proof.* See e.g. [8]. □

Now we may derive the continuity of  $F^p$  as a function of  $\xi$  on  $\text{conv}(\Gamma)$ .

**Theorem 1.** *Let  $\Gamma = \{x_1, \dots, x_k\} \subset \mathbb{R}^d$ ,  $k \in \mathbb{N}$  be a grid. Then*

$$f_\Gamma : \text{conv}(\Gamma) \rightarrow \mathbb{R}, \quad \xi \mapsto F^p(\xi; \Gamma)$$

*is continuous for every  $\xi \in \text{conv}(\Gamma)$ .*

*Proof.* The lower semi-continuity (l.s.c.) of  $f_\Gamma$  follows directly from its dual representation

$$f_\Gamma(\xi) = \max_{u \in \mathbb{R}^d} \min_{1 \leq i \leq k} \{ \|\xi - x_i\|^p + u^T(\xi - x_i) \}$$

and the fact that the maximum of a family of l.s.c. functions is always l.s.c..

To prove that  $f_\Gamma$  is also upper semi-continuous, let  $\xi, \xi^n \in \text{conv}(\Gamma)$ ,  $n \in \mathbb{N}$ , such that  $\xi^n \rightarrow \xi$  as  $n \rightarrow \infty$ . We want to show

$$\limsup_{n \rightarrow \infty} f_\Gamma(\xi^n) \leq f_\Gamma(\xi).$$

Since  $\xi, \xi^n \in \text{conv}(\Gamma)$ , we know that  $f_\Gamma(\xi)$  and  $\limsup_{n \rightarrow \infty} f_\Gamma(\xi^n)$  are finite. Moreover, there is an optimal basis  $I^* \in \mathcal{I}(\Gamma)$ , such that

$$\lambda_{I^*}^* = A_{I^*}^{-1} \begin{bmatrix} \xi \\ 1 \end{bmatrix} \geq 0 \quad \text{and} \quad f_\Gamma(\xi) = \sum_{j \in I^*} \lambda_j^* \|\xi - x_j\|^p.$$

Assume now  $\lambda_{I^*}^* > 0$ . Then there exists some  $n_0 \in \mathbb{N}$  such that for every  $n \geq n_0$ ,

$$\lambda_I^n = A_I^{-1} \begin{bmatrix} \xi^n \\ 1 \end{bmatrix} > 0 \quad \forall n \geq n_0.$$

Hence, we obtain that, for  $n \geq n_0$ ,  $\lambda^n$  (after filling up with zeros) is feasible for (LP) so that

$$\forall n \geq n_0 \quad f_\Gamma(\xi^n) \leq \sum_{j \in I^*} \lambda_j^n \|\xi^n - x_j\|^p.$$

Since  $\lambda_I^n \rightarrow \lambda_I^*$  as  $n \rightarrow \infty$ , this implies

$$\limsup_{n \rightarrow \infty} f_\Gamma(\xi^n) \leq \limsup_{n \rightarrow \infty} \sum_{j \in I^*} \lambda_j^n \|\xi^n - x_j\|^p = f_\Gamma(\xi).$$

In the case  $\lambda_{I^*}^* \not> 0$ , the (LP) is degenerated and we have  $|J^*| < d + 1$  for

$$J^* = \{j \in \{1, \dots, k\} : \lambda_j^* > 0\}$$

as well as  $\xi \in \text{conv}\{x_j : j \in J^*\}$ .

Subsequently, we denote the set of additional bases which are optimal for  $\xi$ , by

$$\mathcal{I}^* = \{I \in \mathcal{I}(\Gamma) : I \supset J^*\},$$

which is a nonempty set by the incomplete basis theorem. It then holds

$$B(\xi, \varepsilon) \cap \text{conv}(\Gamma) \subset \bigcup_{I \in \mathcal{I}^*} \text{conv}\{x_j : j \in I\}$$

for any  $\varepsilon > 0$  small enough.

Let us now fix a subsequence  $n'$  such that

$$f_\Gamma(\xi^{n'}) \rightarrow \limsup_{n \rightarrow \infty} f_\Gamma(\xi^n).$$

Since there are only finite many bases  $I \in \mathcal{I}^*$ , we may choose another subsequence also denoted  $n'$  such that

$$\xi^{n'} \in \text{conv}\{x_j : j \in I'\}$$

for a single optimal basis  $I' \in \mathcal{I}^*$ .

But the last condition yields

$$\lambda_I^{n'} = A_{I'}^{-1} \begin{bmatrix} \xi^{n'} \\ 1 \end{bmatrix} \geq 0,$$

so that again  $\lambda_I^{n'}$  is a feasible solution for the (LP)

$$\begin{aligned} f_\Gamma(\xi^{n'}) &= \min_{\lambda \in \mathbb{R}^k} \sum_{i=1}^k \lambda_i \|\xi^{n'} - x_i\|^p . \\ \text{s.t. } & \begin{bmatrix} x_1 & \dots & x_k \\ 1 & \dots & 1 \end{bmatrix} \lambda = \begin{bmatrix} \xi^{n'} \\ 1 \end{bmatrix}, \lambda \geq 0 \end{aligned}$$

This implies as in the previous case

$$\limsup_{n \rightarrow \infty} f_\Gamma(\xi^n) = \lim_{n' \rightarrow \infty} f_\Gamma(\xi^{n'}) \leq \lim_{n' \rightarrow \infty} \sum_{j \in I'} \lambda_j^{n'} \|\xi^{n'} - x_j\|^p = f_\Gamma(\xi).$$

□

We can now state the main result about the optimality regions  $D_I(\Gamma)$ .

**Proposition 6.** *For  $D_I(\Gamma)$ ,  $I \in \mathcal{I}(\Gamma)$ , it holds*

- (a)  $\{x_j : j \in I\} \subset D_I(\Gamma) \subset \text{conv}\{x_j : j \in I\}$ ,
- (b)  $D_I(\Gamma)$  is closed and therefore a Borel set.

*Proof.* The first assertion follows directly from the definition of  $D_I(\Gamma)$ . To recognize that  $D_I(\Gamma)$  is closed, note that the mappings

$$\xi \rightarrow \sum_{j \in I} \lambda_j^* \|\xi - x_j\|^p \quad \text{and} \quad \xi \rightarrow F^p(\xi; \Gamma)$$

are continuous.

□

Moreover, since (LP) is solvable for every  $\xi \in \text{conv}(\Gamma)$  we clearly have

$$\bigcup_{I \in \mathcal{I}(\Gamma)} D_I(\Gamma) = \text{conv}(\Gamma),$$

*i.e.*  $(D_I(\Gamma))_{I \in \mathcal{I}(\Gamma)}$  provides a Borel measurable covering of  $\text{conv}(\Gamma)$ .

### 2.3 Intrinsic stationarity

To establish the link between the above definition of dual quantization and stationary quantization rules, we have to precise the notion of intrinsic stationarity.

**Definition 3.** (a) *Let  $\Gamma \subset \mathbb{R}^d$ ,  $|\Gamma| < \infty$  be a finite subset of  $\mathbb{R}^d$  and let  $(\Omega_0, \mathcal{S}_0, \mathbb{P}_0)$  be a probability space. Any random operator  $\mathcal{J}_\Gamma : (\Omega_0 \times D, \mathcal{S}_0 \otimes \mathcal{B}^d) \rightarrow \Gamma$ ,  $\text{conv}(\Gamma) \subset D \subset \mathbb{R}^d$  is called a splitting operator (on  $\Gamma$ ).*

*A splitting operator on  $\Gamma$  satisfying*

$$\mathbb{E}_{\mathbb{P}_0}(\mathcal{J}_\Gamma(\xi)) = \int_{\Omega_0} \mathcal{J}_\Gamma(\omega_0, \xi) \mathbb{P}_0(d\xi) = \xi, \quad \forall \xi \in \text{conv}(\Gamma),$$

*is called an intrinsic stationary splitting operator.*

We will see in the next paragraph that  $(\Omega_0, \mathcal{S}_0, \mathbb{P}_0)$  can be modelled as an exogenous probability space in order to randomly “split” (e.g. by simulation) a r.v.  $X$ , defined on the probability space of interest  $(\Omega, \mathcal{S}, \mathbb{P})$ , between the points in  $\Gamma$ .

This new stationarity property is in fact equivalent to the dual stationarity property (9) on the product space  $(\Omega_0 \times \Omega, \mathcal{S}_0 \otimes \mathcal{S}, \mathbb{P}_0 \otimes \mathbb{P})$  as emphasized by the following easy proposition.

**Proposition 7.** *Let  $\text{conv}(\Gamma) \subset D \subset \mathbb{R}^d$ . A random splitting operator  $\mathcal{J}_\Gamma : \Omega_0 \times D \rightarrow \Gamma$  is intrinsic stationary, iff, for any r.v.  $Y : (\Omega, \mathcal{S}, \mathbb{P}) \rightarrow (\mathbb{R}^d, \mathcal{B}^d)$  satisfying  $\text{supp}(\mathbb{P}_Y) \subset \text{conv}(\Gamma)$ ,*

$$\mathbb{E}_{\mathbb{P}_0 \otimes \mathbb{P}}(\mathcal{J}_\Gamma(Y)|Y) = Y \quad \mathbb{P}_0 \otimes \mathbb{P}\text{-a.s.} \quad (11)$$

where  $\mathcal{J}_\Gamma$  and  $Y$  are canonically extended onto  $\Omega_0 \times \Omega$  by setting  $\mathcal{J}_\Gamma((\omega_0, \omega), \cdot) = \mathcal{J}_\Gamma(\omega_0, \cdot)$  and  $Y(\omega_0, \omega) = Y(\omega)$ .

*Proof.* The first implication follows directly from Fubini’s theorem and for the second one simply set  $Y \equiv \xi$ .  $\square$

### 2.3.1 The dual quantization operator $\mathcal{J}_\Gamma^*$

A way to define such an intrinsic stationary random operator in an optimal manner is given by the dual quantization operator  $\mathcal{J}_\Gamma^*$ .

Therefore, let  $\Gamma = \{x_1, \dots, x_k\} \subset \mathbb{R}^d$ ,  $k \in \mathbb{N}$  and assume that  $\text{aff. dim}\{\Gamma\} = d$ . Otherwise the dual quantization operator is not defined.

We then may choose a Borel partition  $(C_I(\Gamma))_{I \in \mathcal{I}(\Gamma)}$  of  $\text{conv}\{\Gamma\}$  such that for every  $I \in \mathcal{I}(\Gamma)$

$$C_I(\Gamma) \subset D_I(\Gamma) = \left\{ \xi \in \mathbb{R}^d : \lambda_I^* := A_I^{-1} \begin{pmatrix} \xi \\ 1 \end{pmatrix} \geq 0 \text{ and } \sum_{j \in I} \lambda_j^* \|\xi - x_j\|^p = F^p(\xi; \Gamma) \right\}.$$

As a consequence and after a reordering of rows,

$$\lambda^I(\xi) = \begin{bmatrix} A_I^{-1} \begin{bmatrix} \xi \\ 1 \end{bmatrix} \\ 0 \end{bmatrix}$$

gives a optimal solution to  $F^p(\xi; \Gamma)$  for every  $\xi \in C_I$ .

Now we are in position to define the intrinsic stationary splitting operator.

**Definition 4.** *Let  $(\Omega_0, \mathcal{S}_0, \mathbb{P}_0) = ([0, 1], \mathfrak{B}([0, 1]), \lambda^1)$  and  $U = \text{Id}_{[0,1]} \sim \mathcal{U}([0, 1])$ . The dual quantization operator  $\mathcal{J}_\Gamma^* : \Omega_0 \times \text{conv}(\Gamma) \rightarrow \Gamma$  is then defined as*

$$\mathcal{J}_\Gamma^*(\omega_0, \xi) = \sum_{I \in \mathcal{I}(\Gamma)} \left[ \sum_{i=1}^k x_i \cdot \mathbb{1}_{\left\{ \sum_{j=1}^{i-1} \lambda_j^I(\xi) \leq U(\omega_0) < \sum_{j=1}^i \lambda_j^I(\xi) \right\}} \right] \mathbb{1}_{C_I(\Gamma)}(\xi). \quad (12)$$

Since

$$\mathbb{E}_{\mathbb{P}_0} \left( \mathbb{1}_{\left\{ \sum_{j=1}^{i-1} \lambda_j^I(\xi) \leq U < \sum_{j=1}^i \lambda_j^I(\xi) \right\}} \right) = \lambda_i^I(\xi)$$

and

$$\forall \xi \in C_I(\Gamma), \quad \sum_{i=1}^k \lambda_i^I(\xi) x_i = \xi,$$

its is clear that  $\mathcal{J}_\Gamma^*$  shares the intrinsic stationarity property:

$$\forall \xi \in \text{conv}\{\Gamma\}, \quad \mathbb{E}_{\mathbb{P}_0}(\mathcal{J}_\Gamma^*(\xi)) = \sum_{I \in \mathcal{I}(\Gamma)} \left[ \sum_{i=1}^k \lambda_i^I(\xi) x_i \right] \mathbb{1}_{C_I(\Gamma)}(\xi) = \xi.$$

On the other hand, one easily checks, that this construction also yields

$$\forall \xi \in \text{conv}\{\Gamma\}, \quad \mathbb{E}_{\mathbb{P}_0} \|\xi - \mathcal{J}_\Gamma^*(\xi)\|^p = F^p(\xi; \Gamma). \quad (13)$$

CHANGE OF NOTATION. From now on, we pass to the product space  $(\Omega_0 \times \Omega, \mathcal{S}_0 \otimes \mathcal{S}, \mathbb{P}_0 \otimes \mathbb{P})$ , where we will also use, if no ambiguity, the symbols  $\mathbb{P}$  and  $\mathbb{E}$  to denote the probability and the expectation on  $\Omega_0 \times \Omega$ . We therefore may assume that the intrinsic stationary splitting operator is independent of all “endogenous” r.v. of interest originally defined on  $(\Omega, \mathcal{S}, \mathbb{P})$  and extended to  $(\Omega_0 \times \Omega, \mathcal{S}_0 \otimes \mathcal{S}, \mathbb{P}_0 \otimes \mathbb{P})$  (which implies that the stationary property (11) holds).

### 2.3.2 Characterizations of the optimal dual quantization error

We use this operator to prove the analogous theorem for dual quantization to Proposition 1.

**Theorem 2.** *Let  $X \in L^0(\Omega, \mathcal{S}, \mathbb{P})$  and  $n \in \mathbb{N}$ . Then*

$$\begin{aligned} d_{n,p}(X) &= \inf \{ \mathbb{E} \|X - \mathcal{J}_\Gamma(X)\|_p : \mathcal{J}_\Gamma : \Omega_0 \times \mathbb{R}^d \rightarrow \Gamma, \text{ intrinsic stationary,} \\ &\quad \text{supp}(\mathbb{P}_X) \subset \text{conv}(\Gamma), |\Gamma| \leq n \} \\ &= \inf \{ \mathbb{E} \|X - \widehat{Y}\|_p : \widehat{Y} : (\Omega_0 \times \Omega, \mathcal{S}_0 \otimes \mathcal{S}, \mathbb{P}_0 \otimes \mathbb{P}) \rightarrow \mathbb{R}^d, \\ &\quad |\widehat{Y}(\Omega_0 \times \Omega)| \leq n, \mathbb{E}(\widehat{Y}|X) = X \} \leq +\infty. \end{aligned}$$

These quantities are finite iff  $X \in L^\infty(\Omega, \mathcal{S}, \mathbb{P})$  and  $n \geq d + 1$ .

*Proof.* First we show the inequality

$$d_n^p(X) \geq \inf \{ \mathbb{E} \|X - \mathcal{J}_\Gamma(X)\|^p : \mathcal{J}_\Gamma : \mathbb{R}^d \rightarrow \Gamma \text{ is intrinsic stationary,} \quad (14)$$

$$\text{supp}(\mathbb{P}_X) \subset \text{conv}(\Gamma), |\Gamma| \leq n \}.$$

We may assume that  $d_n^p(X) < +\infty$  which implies the existence of a grid  $\Gamma \in \mathbb{R}^d$  with  $|\Gamma| \leq n$  and  $d^p(X; \Gamma) < +\infty$  so that Proposition 3 implies  $\text{supp}(\mathbb{P}_X) \subset \text{conv}\{\Gamma\}$ .

Hence, we choose a Borel partition  $(C_I(\Gamma))_{I \in \mathcal{I}(\Gamma)}$  of  $\text{conv}\{\Gamma\}$  with  $C_I(\Gamma) \subset D_I(\Gamma)$ ,  $I \in \mathcal{I}(\Gamma)$ , so that the dual quantization operator  $\mathcal{J}_\Gamma^*$  is well defined by (12).

Owing to the independence of  $X$  and  $\mathcal{J}_\Gamma^*$  on  $\Omega_0 \times \Omega$ , it holds

$$\mathbb{E}(\|\xi - \mathcal{J}_\Gamma^*(\xi)\|^p)_{|\xi=X} = \mathbb{E}(\|X - \mathcal{J}_\Gamma^*(X)\|^p | X) \quad a.s.,$$

so that we conclude from (13)

$$\begin{aligned} \mathbb{E} F^p(X; \Gamma) &= \mathbb{E}[\mathbb{E}(F^p(X; \Gamma) | X)] = \mathbb{E}[\mathbb{E}(F^p(\xi; \Gamma))_{|\xi=X}] \\ &= \mathbb{E}[\mathbb{E}(\|\xi - \mathcal{J}_\Gamma^*(\xi)\|^p)_{|\xi=X}] = \mathbb{E}[\mathbb{E}(\|X - \mathcal{J}_\Gamma^*(X)\|^p | X)] \\ &= \mathbb{E}\|X - \mathcal{J}_\Gamma^*(X)\|^p. \end{aligned}$$

Since  $\mathcal{J}_\Gamma^*$  is intrinsic stationary by construction, the first inequality (14) holds.

The second inequality

$$\begin{aligned} &\inf \{ \mathbb{E} \|X - \mathcal{J}_\Gamma(X)\|^p : \mathcal{J}_\Gamma \text{ is intrinsic stationary, } \text{supp}(\mathbb{P}_X) \subset \text{conv}(\Gamma), |\Gamma| \leq n \} \\ &\geq \inf \{ \mathbb{E} \|X - \widehat{Y}\|^p : \widehat{Y} \text{ is a r.v., } |\widehat{Y}(\Omega_0 \times \Omega)| \leq n, \mathbb{E}(\widehat{Y}|X) = X \}. \end{aligned}$$

follows directly from setting  $\widehat{Y} = \mathcal{J}_\Gamma^*(X)$  in the case  $\mathcal{J}_\Gamma^*$  exists and  $\text{supp}(\mathbb{P}_X) \subset \text{conv}(\Gamma)$ . Otherwise, there is nothing to show.

To prove the reverse inequality, we may assume that there is a r.v.  $\widehat{Y}$  satisfying  $|\widehat{Y}(\Omega_0 \times \Omega)| \leq n$  and

$$\mathbb{E}(\widehat{Y} | X) = X \quad a.s.$$

Denote  $\widehat{Y}(\Omega_0 \times \Omega) = \{y_1, \dots, y_k\}$ ,  $k \leq n$  and let

$$\lambda_i = \mathbb{P}(\widehat{Y} = y_i | X), \quad 1 \leq i \leq k$$

be arbitrary versions of the conditional probabilities.

Hence, there exists a null set  $N \in \mathcal{S}_0 \otimes \mathcal{S}$  such that

$$\forall \bar{\omega} = (\omega_0, \omega) \in N^c : \begin{cases} \sum_{i=1}^k y_i \lambda_i(\bar{\omega}) = \mathbb{E}(\widehat{Y} | X)(\bar{\omega}) = X(\omega) \\ \sum_{i=1}^k \lambda_i(\bar{\omega}) = 1 \\ \lambda_i(\bar{\omega}) \in [0, 1], \quad 1 \leq i \leq k. \end{cases}$$

Setting  $\Gamma = \{x_1, \dots, x_k\}$ , we get for every  $\bar{\omega} \in N^c$

$$\begin{aligned} \mathbb{E}(\|X - \widehat{Y}\|^p | X)(\bar{\omega}) &= \sum_{i=1}^k \lambda_i(\bar{\omega}) \mathbb{E}(\|X - y_i\|^p | X)(\bar{\omega}) = \sum_{i=1}^k \lambda_i(\bar{\omega}) \|X(\omega) - y_i\|^p \\ &\geq F^p(X(\omega); \Gamma). \end{aligned}$$

Taking the expectation completes the proof.  $\square$

*Remark.* We necessarily need to define  $\widehat{Y}$  on the larger product probability space  $(\Omega_0 \times \Omega, \mathcal{S}_0 \otimes \mathcal{S}, \mathbb{P}_0 \otimes \mathbb{P})$  rather than only on  $(\Omega, \mathcal{S}, \mathbb{P})$ , since  $\mathcal{S}$  might not be fine enough to contain appropriated r.v.s  $\widehat{Y}$  satisfying  $\mathbb{E}(\widehat{Y} | X) = X$ . E.g., if  $\mathcal{S} = \sigma(X)$ ,  $\widehat{Y}$  would be  $\sigma(X)$ -measurable so that  $\mathbb{E}(\widehat{Y} | X) = \widehat{Y}$ , intrinsic stationarity would become unreachable for general finite-valued r.v.  $\widehat{Y}$ .

### 2.3.3 Applications of the intrinsic stationarity

As a consequence of the above Theorem we get the following theorem about cubature by dual quantization.

**Theorem 3.** *Let  $X \in L^2(\mathbb{P})$  and assume that  $F \in C^{1,1}(\mathbb{R}^d)$  is differentiable with Lipschitz derivative. For any grid  $\Gamma = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$  with  $\text{conv}\{\Gamma\} \supset \text{supp}(\mathbb{P}_X)$  it holds for the cubature formula  $\mathbb{E}F(\mathcal{J}_\Gamma^*(X)) = \sum_{i=1}^n \mathbb{P}(\mathcal{J}_\Gamma^*(X) = x_i) \cdot F(x_i)$*

$$|\mathbb{E}F(X) - \mathbb{E}F(\mathcal{J}_\Gamma^*(X))| \leq [F']_{Lip} \mathbb{E}\|X - \mathcal{J}_\Gamma^*(X)\|^2.$$

Now assume that the integrand  $F$  is convex. if  $\widehat{X}^\Gamma$  is a quantization which satisfies the regular stationarity property  $\mathbb{E}(X | \widehat{X}^\Gamma) = \widehat{X}^\Gamma$ , it follows from Jensen's inequality that  $\mathbb{E}F(\widehat{X}^\Gamma)$  yields a lower bound for the approximation of  $\mathbb{E}F(X)$ .

In contrast to that and exploiting the intrinsic stationarity of  $\mathcal{J}_\Gamma^*$ , a cubature formula based on  $\mathcal{J}_\Gamma^*$  yields for convex  $F$  an upper bound, which is now valid for any grid  $\Gamma \subset \mathbb{R}^d$ .

**Proposition 8.** *Let  $X \in L^p(\mathbb{P})$ ,  $p \geq 1$  and assume that  $F$  is convex. Then it holds for any grid with  $\Gamma \subset \mathbb{R}^d$  and  $\text{conv}\{\Gamma\} \supset \text{supp}(\mathbb{P}_X)$*

$$\mathbb{E}F(\mathcal{J}_\Gamma^*(X)) \geq \mathbb{E}F(X).$$

*Proof.* It follows from the intrinsic stationarity  $\mathbb{E}(\mathcal{J}_\Gamma^*(X) | X) = X$  and Jensen's inequality for conditional expectations that

$$\mathbb{E}F(X) = \mathbb{E}[F(\mathbb{E}(\mathcal{J}_\Gamma^*(X) | X))] \leq \mathbb{E}[\mathbb{E}(F(\mathcal{J}_\Gamma^*(X)) | X)] = \mathbb{E}F(\mathcal{J}_\Gamma^*(X)).$$

$\square$

## 2.4 Upper bounds and product quantization

**Proposition 9** (Scalar bound). *Let  $\Gamma = \{x_1, \dots, x_n\} \subset \mathbb{R}$  with  $x_1 \leq \dots \leq x_n$ . Then*

$$\forall \xi \in [x_1, x_n], \quad F^p(\xi, \gamma) \leq \max_{1 \leq i \leq n-1} \left( \frac{x_{i+1} - x_i}{2} \right)^p.$$

*Proof.* If  $\xi \in \Gamma$ , then  $F^p(\xi, \gamma) = 0$  and the assertion holds. Suppose now  $\xi \in (x_i, x_{i+1})$ . Then  $\xi = \lambda x_i + (1 - \lambda)x_{i+1}$  and  $\lambda = \frac{x_{i+1} - \xi}{x_{i+1} - x_i}$ , so that

$$F^p(\xi, \Gamma) \leq \left( \frac{x_{i+1} - \xi}{x_{i+1} - x_i} \right) |\xi - x_i|^p + \left( \frac{\xi - x_i}{x_{i+1} - x_i} \right) |\xi - x_{i+1}|^p$$

attains its maximum at  $\xi = \frac{x_i + x_{i+1}}{2}$ . This implies

$$F^p(\xi, \Gamma) \leq \left( \frac{1}{2} + \frac{1}{2} \right) \left| \frac{x_{i+1} - x_i}{2} \right|^p,$$

which yields the assertion.  $\square$

**Proposition 10** (Local product Quantization). *Let  $\|\cdot\| = |\cdot|_p$  be the canonical  $p$ -norm on  $\mathbb{R}^d$ ,  $\xi = (\xi_1, \dots, \xi_d) \in \mathbb{R}^d$  and  $\Gamma = \prod_{j=1}^d \alpha_j$  for some  $\alpha_j \subset \mathbb{R}$ . Then*

$$F^p(\xi; \Gamma) = \sum_{j=1}^d F^p(\xi_j; \alpha_j).$$

*Proof.* Denoting  $\alpha_j = \{a_1^j, \dots, a_{n_j}^j\}$ ,  $\Gamma = \{x_1, \dots, x_n\}$  and due to the fact that  $\{x_1, \dots, x_n\}$  is made up by the cartesian product of  $\{a_1^j, \dots, a_{n_j}^j\}$ ,  $j = 1, \dots, d$  we have for any  $u, \xi \in \mathbb{R}^d$ :

$$\min_{1 \leq i \leq n} \left\{ \sum_{j=1}^d |\xi_j - x_i^j|^p + u_j (\xi_j - x_i^j) \right\} = \sum_{j=1}^d \min_{1 \leq i \leq n_j} \{ |\xi_j - a_i^j|^p + u_j (\xi_j - a_i^j) \}.$$

We then get from Proposition 4

$$\begin{aligned} F^p(\xi; \Gamma) &= \max_{u \in \mathbb{R}^d} \min_{1 \leq i \leq n} \left\{ \sum_{j=1}^d |\xi_j - x_i^j|^p + u_j (\xi_j - x_i^j) \right\} \\ &= \max_{u \in \mathbb{R}^d} \sum_{j=1}^d \min_{1 \leq i \leq n_j} \{ |\xi_j - a_i^j|^p + u_j (\xi_j - a_i^j) \} \\ &= \sum_{j=1}^d \max_{u_j \in \mathbb{R}} \min_{1 \leq i \leq n_j} \{ |\xi_j - a_i^j|^p + u_j (\xi_j - a_i^j) \} \\ &= \sum_{j=1}^d F^p(\xi_j; \alpha_j). \end{aligned}$$

$\square$

This enables us to derive a first upper bound for the asymptotics of the optimal dual quantization error of distributions with bounded support when the size of the grid tends to infinity.

**Proposition 11** (Product Quantization). *Let  $C = a + \ell[0, 1]^d$ ,  $a = (a_1, \dots, a_d) \in \mathbb{R}^d$ ,  $\ell > 0$ , be a hypercube, parallel to the coordinate axis with common edge length  $\ell$ . Let  $\Gamma$  be the product quantizer of size  $(m + 1)^d$  defined by*

$$\Gamma = \prod_{k=1}^d \left\{ a_j + \frac{i\ell}{m}, i = 0, \dots, m \right\}.$$

Then it holds

$$\forall \xi \in C, \quad F_p^p(\xi; \Gamma) \leq d \cdot C_{\|\cdot\|} \cdot \left(\frac{\ell}{2}\right)^p \cdot m^{-p} \quad (15)$$

with some constant  $C_{\|\cdot\|} > 0$ . Moreover, for any compactly supported r.v.  $X$

$$d_{n,p}(X) = \mathcal{O}(n^{-1/d}).$$

*Proof.* The first claim follows directly from Propositions 9 and 10. For the second assertion let  $n \geq 2^d$  and set  $m = \lfloor n^{1/d} \rfloor - 1$ . If we choose the hypercube  $C$  such  $\text{supp}(\mathbb{P}_X) \subset C$  we arrive at

$$d_n^p(X) \leq C_1 \left( \frac{1}{\lfloor n^{1/d} \rfloor - 1} \right)^p \leq C_2 \left( \frac{1}{n} \right)^{p/d}$$

for some constants  $C_1, C_2 > 0$ , which yields the desired upper bound.  $\square$

## 2.5 Extension for distributions with unbounded support

We have seen in the previous sections, that  $F^p(\xi; \Gamma)$  is finite if and only if  $\xi \in \text{conv}\{\Gamma\}$ , so that intrinsic stationarity cannot hold for a r.v.  $X$  with unbounded support.

Nevertheless, we may restrict the stationarity requirement in the definition of the dual quantization error for unbounded  $X$  to its “natural domain”  $\text{conv}(\Gamma)$ , which means that we drop the constraint  $\text{supp}(\mathbb{P}_X) \subset \text{conv}(\Gamma)$  from Theorem 2.

**Definition 5.** *We define the extended dual  $L^p$ -mean quantization error as*

$$\bar{d}_n^p(X) = \inf \{ \mathbb{E} \|X - \mathcal{J}_\Gamma(X)\|^p : \mathcal{J}_\Gamma : \mathbb{R}^d \rightarrow \Gamma \text{ is intrinsic stationary}, \Gamma \subset \mathbb{R}^d, |\Gamma| \leq n \}.$$

Combining Propositions 1 and Theorem 2 we get

**Proposition 12.** *Let  $X \in L^p(\mathbb{P})$ . Then*

$$\bar{d}_n^p(X) = \inf \{ \mathbb{E} \bar{F}^p(X; \Gamma) : \Gamma \subset \mathbb{R}^d, |\Gamma| \leq n \}$$

for

$$\bar{F}^p(\xi; \Gamma) = F^p(\xi; \Gamma) \mathbb{1}_{\text{conv}(\Gamma)}(\xi) + \|\xi - \pi_\Gamma(\xi)\|^p \mathbb{1}_{\text{conv}(\Gamma)^c}(\xi).$$

Note, that we have for any  $X \in L^p(\mathbb{P})$

$$\bar{d}_n^p(X) \leq d_n^p(X),$$

where equality in general even does not hold anymore for  $X$  with bounded support, but it was shown in another paper ([11]), that both quantities coincide asymptotically in the bounded case.



## 2.6 Sharp rate of convergence : Zador's Theorem for dual quantization

In the companion paper [11], we establish the following theorem which looks formally identical to the celebrated Zador Theorem for regular vector quantization.

**Theorem 4.** (a) Let  $X \in L^{p+\delta}(\mathbb{P})$ ,  $\delta > 0$ , absolutely continuous w.r.t. to  $\lambda^d$  and  $\mathbb{P}_X = h\lambda^d$ . Then

$$\lim_{n \rightarrow \infty} n^{1/d} \bar{d}_{n,p}(X) = Q_{d,p,\|\cdot\|} \cdot \|h\|_{d/(d+p)}^{1/p}$$

where

$$Q_{d,p,\|\cdot\|} = \lim_{n \rightarrow \infty} n^{1/d} \bar{d}_{n,p}(\mathcal{U}([0,1]^d)) = \inf_{n \geq 1} n^{1/d} \bar{d}_{n,p}(\mathcal{U}([0,1]^d)).$$

This constant satisfies  $Q_{d,p,\|\cdot\|} \geq Q_{d,p,\|\cdot\|}^{vq}$ , where  $Q_{d,p,\|\cdot\|}^{vq}$  denotes the asymptotic constant for the sharp Voronoi vector quantization rate of the uniform distribution over  $[0,1]^d$ , i.e.

$$Q_{d,p,\|\cdot\|}^{vq} = \lim_{n \rightarrow \infty} n^{1/d} e_{n,p}(\mathcal{U}([0,1]^d)) = \inf_{n \geq 1} n^{1/d} e_{n,p}(\mathcal{U}([0,1]^d)).$$

Furthermore, when  $d = 1$  we know that  $Q_{d,p,\|\cdot\|} = (\frac{2^{p+1}}{p+2})^{1/p} Q_{d,p,\|\cdot\|}^{vq}$ .

(b) When  $X$  has a compact support the above sharp rate holds for  $d_{n,p}(X)$  as well.

## 3 Quadratic Euclidean case and Delaunay Triangulation

In the case that  $(\mathbb{R}^d, \|\cdot\|)$  is the Euclidean space and  $p = 2$ , the optimality regions  $D_I(\Gamma)$  have either empty interior or are maximal, i.e.  $\mathring{D}_I(\Gamma) = \emptyset$  or  $D_I(\Gamma) = \text{conv}\{x_j; j \in I\}$ . This follows from the fact that in the quadratic Euclidean case the dual feasibility of a basis  $I \in \mathcal{I}(\Gamma)$  can be stated independently of  $\xi$ .

This feature is also the key to the following theorem, which was first proved by Rajan in [13] and establishes the link between a solution to  $F^2(\xi; \Gamma)$  (the so-called power function in [13]) and the Delaunay property of a triangle.

Recall that a triangle (or  $d$ -simplex)  $\text{conv}\{x_{i_1}, \dots, x_{i_{d+1}}\}$  in a set of points  $\Gamma = \{x_1, \dots, x_k\}$ ,  $k \geq d+1$  has the *Delaunay property*, if the sphere spanned by  $\{x_{i_1}, \dots, x_{i_{d+1}}\}$  contains no point of  $\Gamma$  in its interior.

**Theorem 5.** Let  $\|\cdot\| = |\cdot|_2$  be the Euclidean norm,  $p = 2$ , and  $\Gamma = \{x_1, \dots, x_k\} \subset \mathbb{R}^d$  with  $\text{aff. dim}\{\Gamma\} = d$ .

(a) If  $I \in \mathcal{I}(\Gamma)$  defines a Delaunay triangle (or  $d$ -simplex), then

$$\lambda_I = A_I^{-1} \begin{pmatrix} \xi \\ 1 \end{pmatrix}$$

provides a solution to LP for every  $\xi \in \text{conv}\{x_j : j \in I\}$ .

In particular, this implies  $D_I(\Gamma) = \text{conv}\{x_j : j \in I\}$ .

(b) If  $I \in \mathcal{I}(\Gamma)$  satisfies  $\mathring{D}_I(\Gamma) \neq \emptyset$ , then the triangle (or  $d$ -simplex) defined by  $I$  has the Delaunay property for  $\Gamma$ .

We provide here a short proof based on the duality for Linear Programs for the reader's convenience.

*Proof.* First note, that  $I \in \mathcal{I}(\Gamma)$  defines a Delaunay triangle, if there is exists a center  $z \in \mathbb{R}^d$  such that for every  $j \in I$

$$|z - x_j|_2 \leq |z - x_i|_2, \quad 1 \leq i \leq k, \quad (16)$$

and equality holds for  $i \in I$ .

Suppose that  $z = \xi + \frac{u_1}{2}$ . Then

$$\forall i \in I, \quad |z - x_i|_2^2 = |\xi - x_i|_2^2 + \xi^T u_1 - x_i^T u_1 + \left| \frac{u_1}{2} \right|_2^2$$

so that (16) is equivalent to

$$\begin{aligned} |\xi - x_j|_2^2 - x_j^T u_1 &\leq |\xi - x_i|_2^2 - x_i^T u_1, \quad 1 \leq i \leq k, j \in I, \\ u_2 &= |\xi - x_j|_2^2 - x_j^T u_1, \quad j \in I. \end{aligned} \tag{17}$$

Note that this exactly the dual feasibility condition of Proposition 5.

(a) Now let  $I \in \mathcal{I}(\Gamma)$  such that  $\{x_j : j \in I\}$  defines a Delaunay triangle. Denoting by  $z \in \mathbb{R}^d$  the center of the sphere spanned by  $\{x_j; j \in I\}$ , we define  $u = (u_1, u_2)$  for every  $\xi \in \mathbb{R}^d$  as

$$u_1 = 2(z - \xi) \quad \text{and} \quad u_2 = |\xi - x_j|_2^2 - x_j^T u_1 \quad \text{for an arbitrary } j \in I.$$

Consequently  $z = \xi + \frac{u_1}{2}$ , so that  $u$  is dual feasible for (LP) due to the above said.

Since  $\lambda_I = A_I^{-1} \begin{pmatrix} \xi \\ 1 \end{pmatrix} \geq 0$  iff  $\xi \in \text{conv}\{x_j : j \in I\}$ , Proposition 5(a) then yields that  $\lambda_I$  provides an optimal solution to (LP) for any  $\xi \in \text{conv}\{x_j : j \in I\}$ .

(b) Let  $I \in \mathcal{I}(\Gamma)$  and choose some  $\xi \in \overset{\circ}{D}_I(\Gamma)$ . Then Proposition 6(a) implies  $\xi \in \overbrace{\text{conv}\{x_j : j \in I\}}^{\circ}$ . As a consequence, it holds  $\lambda_I = A_I^{-1} \begin{pmatrix} \xi \\ 1 \end{pmatrix} > 0$ , so that we conclude from Proposition 5(b) that the unique dual solution to (LP) is given by  $\begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = (A_I^T)^{-1} c_I$ . Since moreover  $A^T \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} \leq c$ ,  $(u_1, u_2)$  satisfies (17) so that

$$z = \xi + \frac{u_1}{2}$$

is the center of a Delaunay triangle containing  $\xi$  in its interior. □

Consequently, if a grid  $\Gamma \subset \mathbb{R}^d$  exhibits a Delaunay triangulation, the dual quantization operator  $\mathcal{J}_\Gamma^*$  is (up to the triangles borders) uniquely defined and maps any  $\xi \in \text{conv}(\Gamma)$  to the vertices of the Delaunay triangle in which  $\xi$  lies.

This yields a duality relation of  $\mathcal{J}_\Gamma^*$  and the nearest neighbor projection  $\pi_\Gamma$ , which is based on the Voronoi tessellation - the dual counterpart of the Delaunay triangulation in the graph theoretic sense.

## 4 Existence of an optimal dual quantization grid

In order to derive the existence of the optimal dual quantization grids, i.e. the fact that the infimum over all grids  $\Gamma \subset \mathbb{R}^d$  with  $|\Gamma| \leq n$  in Definition 2 holds actually as a minimum, we have to discuss properties of  $F_p$  and  $d_p$  as mapping of the quantization grid  $\Gamma$ .

We therefore define for every  $n \geq 1$  and every  $\gamma = (x_1, \dots, x_n) \in (\mathbb{R}^d)^n$

$$F_{n,p}(\xi, \gamma) = \inf \left\{ \left( \sum_{1 \leq i \leq n} \lambda_i \|\xi - x_i\|^p \right)^{1/p} : \lambda_i \in [0, 1] \text{ and } \sum_{1 \leq i \leq n} \lambda_i x_i = \xi, \sum_{1 \leq i \leq n} \lambda_i = 1 \right\}$$

and

$$d_{n,p}(X, \gamma) = \|F_{n,p}(X, \gamma)\|_{L^p}.$$

This functions are clearly symmetric and in fact *only depend on the value set* of  $\gamma = (x_1, \dots, x_n)$  denoted  $\Gamma = \Gamma_\gamma = \{x_i, i = 1, \dots, n\}$ . Hence, we have

$$F_{n,p}(\xi, \gamma) = F_p(\xi; \Gamma_\gamma) \quad \text{and} \quad d_{n,p}(X, \gamma) = d_p(X; \Gamma_\gamma),$$

which implies

$$d_{n,p}(X) = \inf \{d_{n,p}(X, \gamma) : \gamma \in (\mathbb{R}^d)^n\}.$$

One also carries over these definitions to the unbounded case, i.e. we obtain  $\bar{F}_{p,n}(\xi, \gamma)$  and  $\bar{d}_{n,p}(X, \gamma)$ .

As in section 2, we may drop a duplicate parameter  $p$  in the  $p$ -th power of the above expression, e.g. we write  $F_n^p(\xi, \gamma)$  instead of  $F_{n,p}^p(\xi, \gamma)$ . Moreover, we assume again without loss of generality that  $\text{Conv}(\text{supp } \mathbb{P}_X)$  has a nonempty interior in  $\mathbb{R}^d$  or equivalently that

$$\text{span}(\text{supp } \mathbb{P}_X) = \mathbb{R}^d.$$

#### 4.1 Distributions with compact support

We first handle the case when  $\text{supp}(\mathbb{P}_X)$  is compact.

**Theorem 6.** (a) Let  $p \in [1, +\infty)$ . For every integer  $n \geq 1$ , the  $L^p$ -mean dual quantization error function  $\gamma \mapsto d_{n,p}(X, \gamma)$  is l.s.c. and if  $p > 1$  it also attains a minimum.

(b) Let  $p > 1$ . If  $|\text{supp}(\mathbb{P}_X)| \geq n$ , any optimal grid  $\Gamma^{n,*}$  has size  $n$  and  $d_{n,p}(X) = 0$  if and only if  $|\text{supp}(\mathbb{P}_X)| \leq n$ . Furthermore, the sequence  $n \mapsto d_{n,p}(X)$  decreases (strictly) to 0 as long as it does not vanish.

*Proof.* (a) *Lower semi-continuity.* Let  $\gamma^{(k)}$ ,  $k \geq 1$  be a sequence of  $n$ -tuples that converges towards  $\gamma^{(\infty)}$ . The dual representation

$$F_n^p(\xi, (x_1, \dots, x_n)) = \max_{u \in \mathbb{R}^d} \min_{1 \leq i \leq n} \{ \|\xi - x_i\|^p + u^T(\xi - x_i) \}$$

then yields

$$\liminf_{k \rightarrow \infty} F_n^p(\xi, \gamma^{(k)}) \geq F_n^p(\xi, \gamma^{(\infty)}).$$

Consequently, one derives the lower semi-continuity of  $d_{n,p}(X, \cdot)$  since

$$\liminf_k d_n^p(X, \gamma^{(k)}) \geq \mathbb{E} \left( \liminf_k F_n^p(X, \gamma^{(k)}) \right) \geq \mathbb{E} \left( F_n^p(X, \gamma^{(\infty)}) \right) = d_n^p(X, \gamma^{(\infty)})$$

owing to Fatou's lemma.

*Existence of an optimal dual quantization grid.*

Assume that  $\gamma^{(k)} = (x_1^{(k)}, \dots, x_n^{(k)})$ ,  $k \geq 1$ , is a general sequence of  $n$ -tuples such that  $\liminf_k d_{n,p}(X, \gamma^{(k)}) < +\infty$ . Then  $\liminf_k \min_{1 \leq i \leq n} |x_i^{(k)}| < +\infty$  since, otherwise

$$\liminf_{k \rightarrow \infty} d_n^p(X, \gamma^{(k)}) \geq \mathbb{E} \text{dist}(X, \gamma^{(k)})^p \geq \mathbb{E} \liminf_{k \rightarrow \infty} \text{dist}(X, \gamma^{(k)})^p = +\infty$$

owing to Fatou's lemma.

Now, up to appropriate extractions, one may assume that  $d_{n,p}(X, \gamma^{(k)})$  converges to a finite limit and that there exists a nonempty set of indices  $J_\infty$  such that

$$\forall j \in J_\infty, x_j^{(k)} \rightarrow x_j^{(\infty)}, \quad \forall j \notin J_\infty, \|x_j^{(k)}\| \rightarrow +\infty \text{ as } k \rightarrow \infty.$$

Let  $\xi \in \text{supp}(\mathbb{P}_X)$ ,  $\gamma^{(\infty)}$  be any  $n$ -tuple of  $(\mathbb{R}^d)^n$  such that  $\Gamma_{\gamma^{(\infty)}} = \{x_j^{(\infty)}, j \in J_\infty\}$  and denote  $n_\infty = |J_\infty|$ . We then want to show

$$\liminf_{k \rightarrow \infty} F_n^p(\xi, \gamma^{(k)}) \geq F_n^p(\xi, \gamma^{(\infty)}). \quad (18)$$

Moreover, let  $u \in \mathbb{R}^d$  and  $(y_k)$  be a sequence such that  $\|y_k\| \rightarrow +\infty$ . Then it holds for  $p > 1$

$$\|\xi - y_k\|^p + u^T(\xi - y_k) \rightarrow +\infty \quad \text{as } k \rightarrow \infty.$$

In the case when  $u^T(\xi - y_k)$  is bounded from below, the claim is trivial. Otherwise, we have  $u^T(\xi - y_k) \rightarrow -\infty$  so that for  $k$  large enough it holds

$$\|\xi - y_k\|^p + u^T(\xi - y_k) = \|\xi - y_k\|^p - |u^T(\xi - y_k)|.$$

Applying Cauchy-Schwarz and using the equivalence of norms on  $\mathbb{R}^d$  we arrive at

$$\|\xi - y_k\|^p + u^T(\xi - y_k) \geq \|\xi - y_k\|^p - |u|_2 \|\xi - y_k\|_2 \geq \|\xi - y_k\| (\|\xi - y_k\|^{p-1} - C_{\|\cdot\|} \|u\|_2) \rightarrow +\infty.$$

This yields for any  $u \in \mathbb{R}^d$

$$\liminf_{k \rightarrow \infty} \min_{1 \leq i \leq n} \{\|\xi - x_i^{(k)}\|^p + u^T(\xi - x_i^{(k)})\} \geq \min_{i \in J_\infty} \{\|\xi - x_j^{(\infty)}\|^p + u^T(\xi - x_j^{(\infty)})\},$$

so that the dual representation of  $F_n^p$  finally implies (18).

Now, assume that the sequence  $(\gamma^{(k)})_{k \geq 1}$  is asymptotically optimal. Then  $d_{n,p}(X) = \lim_k d_{n,p}(X, \gamma^{(k)}) < +\infty$ . Applying what precedes yields

$$d_{n,p}(X) = \lim_k d_{n,p}(X, \gamma^{(k)}) \geq d_{n_\infty,p}(X, \Gamma_{\gamma^{(\infty)}}) \geq d_{n_\infty,p}(X) \geq d_{n,p}(X)$$

so that

$$d_{n,p}(X) = d_{n_\infty,p}(X, \Gamma_{\gamma^{(\infty)}}) = d_{n_\infty,p}(X).$$

This proves the existence of an optimal dual quantizer at level  $n$ .

(b) To prove that the  $L^p$ -mean dual quantization error decreases with optimal grids of full size  $n$  at level  $n$ , as long as it does not vanish, we will proceed by induction.

CASE  $n = d + 1$ . Then  $J_\infty^c = \emptyset$  and furthermore  $\Gamma_{\gamma^{(\infty)}}$  has size  $d + 1$  since its convex hull contains  $\text{supp}(\mathbb{P}_X)$  which has a nonempty interior. Owing to the lower semi-continuity of the function  $d_{n,p}(X, \cdot)$ ,  $\gamma^{(\infty)}$  is optimal. Furthermore, if  $\text{supp}(\mathbb{P}_X) = \Gamma_{n_0} := \{x_1, \dots, x_{n_0}\}$  has size  $n_0 \leq d + 1$ , then setting for every  $\xi = x_{i_0}$ ,  $\lambda_j = \delta_{i_0,j}$  yields  $F_{n_0,p}(\xi; \Gamma_{n_0}) = 0$ , which implies  $d_{n_0,p}(X) = d_{n_0,p}(X; \Gamma_{n_0}) = 0$ .

CASE  $n > d + 1$ . Assume now that  $|\text{supp}(\mathbb{P}_X)| \geq n$ . Then there exists by the induction assumption an optimal grid  $\Gamma_{n-1}^* = \{x_1^*, \dots, x_{n-1}^*\} \subset \mathbb{R}^d$  at level  $n - 1$  which is optimal for  $d_{n-1,p}(X, \cdot)$  and contains exactly  $n - 1$  points. This grid contains  $d + 1$  affinely independent points since  $d_{n-1,p}(X) < +\infty$  (since  $\text{span}(\text{supp}(\mathbb{P}_X)) = \mathbb{R}^d$ ). Let  $\xi_0 \in \text{supp}(\mathbb{P}_X) \setminus \Gamma_{n-1}^*$  and let  $\Gamma_{n-1}(\xi_0) = \{x_i^*, i \in I_0\}$  be some affinely independent points from  $\Gamma_{n-1}^*$ , solution to the optimization problem (LP) at level  $n - 1$  for  $F_{n-1,p}(\xi_0, \Gamma_{n-1}^*)$ . There exists  $I \subset \{1, \dots, n - 1\}$  such that and

$$I \supset I_0, |I| = d + 1, \{x_i^*, i \in I\} \text{ is an affine basis of } \mathbb{R}^d.$$

By the (affine) exchange lemma, for every index  $j \in I_0$ ,  $\{x_i^*, i \in I, i \neq j\} \cup \{\xi_0\}$  is an affine basis. Furthermore  $\bigcup_{j \in I_0} \left( B(\xi_0; \varepsilon) \cap \text{Conv}(\{x_i^*, i \in I, i \neq j\} \cup \{\xi_0\}) \right)$  is a neighbourhood of  $\xi_0$  in

$\text{Conv}\{\Gamma_{n-1}^*\}$ . Consequently there exists  $i_0 \in I_0$  such that

$$\mathbb{P}\left(X \in B(\xi_0; \varepsilon) \cap \text{Conv}(\{x_i^*, i \in I, i \neq j\} \cup \{\xi_0\})\right) > 0$$

since  $\xi_0 \in \text{supp}(\mathbb{P}_X) \subset \text{Conv}\{\Gamma_{n-1}^*\}$ .

Now for every  $v \in B_{\|\cdot\|}(0; 1)$ ,  $v$  writes on the vector basis  $\{x_i^* - \xi_0\}_{i \in I \setminus \{i_0\}}$ ,  $v = \sum_{i \in I \setminus \{i_0\}} \theta_i (x_i^* - \xi_0)$  with coordinates  $\theta_i$  satisfying  $\sum_{i \in I \setminus \{i_0\}} |\theta_i| \leq C_{d, \|\cdot\|}$ , where  $C_{d, \|\cdot\|} \in [1, +\infty)$  is a real constant only depending on  $d$  and the norm  $\|\cdot\|$ .

Let  $\varepsilon \in (0, \frac{1}{C_{d, \|\cdot\|}})$  be a positive real number to be specified later on.

Let  $\zeta \in B_{\|\cdot\|}(\xi_0; \varepsilon) \cap \text{Conv}(\{x_i^*, i \in I, i \neq i_0\} \cup \{\xi_0\})$ . Then  $v = \frac{\zeta - \xi_0}{\varepsilon} \in B_{\|\cdot\|}(0; 1)$  and

$$\zeta = (1 - \varepsilon \underbrace{\sum_{i \in I \setminus \{i_0\}} \theta_i}_{>0}) \xi_0 + \varepsilon \sum_{i \in I \setminus \{i_0\}} \theta_i x_i^*.$$

Furthermore, by the uniqueness of the decomposition (with sum equal to 1), we also know that  $\theta_i \geq 0$ ,  $i \in I \setminus \{i_0\}$ . Consequently

$$F_n^p(\zeta, \Gamma_{n-1}^* \cup \{\xi\}) \leq (1 - \varepsilon \sum_{i \in I \setminus \{i_0\}} \theta_i) |\zeta - \xi_0|^p + \varepsilon \sum_{i \in I \setminus \{i_0\}} \theta_i |\zeta - x_i^*|^p.$$

Now set  $L^* := \max_{i \in I} |\xi_0 - x_i^*|$ . Then

$$|\zeta - \xi_0| \leq \varepsilon \sum_{i \in I \setminus \{i_0\}} \theta_i |x_i^* - \xi_0| \leq \varepsilon C_{d, \|\cdot\|} L^*$$

and, for every  $i \in I \setminus \{i_0\}$ ,

$$|\zeta - x_i^*| \leq |\zeta - \xi_0| + L^* \leq C_{d, \|\cdot\|} \varepsilon + L^*$$

Finally, for every  $\varepsilon \in (0, \frac{1}{C_{d, \|\cdot\|}})$  and every  $\zeta \in B_{\|\cdot\|}(\xi_0; \varepsilon)$ ,

$$F_n^p(\zeta, \Gamma_{n-1}^* \cup \{\xi_0\}) \leq \varepsilon^p \tilde{L}^*.$$

On the other hand, if  $\varepsilon < d(\xi_0, \Gamma_{n-1}^*)$ ,

$$F_{n-1}^p(\zeta, \Gamma_{n-1}^*) \geq \text{dist}(\zeta, \Gamma_{n-1}^*)^p \geq (\text{dist}(\xi_0, \Gamma_{n-1}^*) - \varepsilon)^p$$

so that, for small enough  $\varepsilon$ ,  $\varepsilon^p \tilde{L}^* < F_{n-1}^p(\zeta, \Gamma_{n-1}^*)$  which finally proves the existence of an  $\varepsilon_0 > 0$  such that

$$\forall \zeta \in B_{\|\cdot\|}(\xi_0; \varepsilon) \cap \text{Conv}(\{x_i^*, i \in I, i \neq i_0\} \cup \{\xi_0\}), \quad F_n^p(\zeta, \Gamma_{n-1}^* \cup \{\xi_0\}) < F_{n-1}^p(\zeta, \Gamma_{n-1}^*).$$

As a first result,

$$d_{n,p}(X) \leq d_p(X; \Gamma_{n-1}^* \cup \{\xi_0\}) < d_p(X; \Gamma_{n-1}^*) = d_{n-1,p}(X).$$

Furthermore, this shows that  ${}^c J_\infty$  is empty *i.e.* all the components of  $\gamma^{(k')}$  remain bounded and converge towards  $\gamma^{(\infty)}$ . Hence  $\gamma^{(\infty)}$  has  $n$  pairwise distinct components since  $d_{n,p}(X; \gamma^{(\infty)}) = d_{n,p}(X) < d_{n-1,p}(X)$  owing to the *l.s.c.*

(c) Convergence to 0 : this follows from Proposition 11. □

## 4.2 Distributions with unbounded support

Let  $X \in L^p(\mathbb{P})$  and let  $r \geq 1$ . We define

$$\bar{F}_p(\xi; \Gamma) = F_p(\xi; \Gamma) \mathbf{1}_{\{X \in \text{Conv}(\Gamma)\}} + \text{dist}(\xi, \Gamma) \mathbf{1}_{\{X \notin \text{Conv}(\Gamma)\}}$$

and

$$\bar{d}_p(X; \Gamma) = \|\bar{F}_p(X; \Gamma)\|_{L^p} < +\infty,$$

since  $d_p(X; \Gamma) \leq \text{diam}(\Gamma) + \|\text{dist}(X, \Gamma)\|_{L^p}$ .

**Theorem 7.** Let  $p > 1$ . Assume that the distribution  $\mathbb{P}_X$  is strongly continuous in the sense that

$$\forall H \text{ hyperplane of } \mathbb{R}^d, \mathbb{P}(X \in H) = 0$$

and has a support with a nonempty interior. Then the extended  $L^p$ -mean dual quantization error function  $\gamma \mapsto \bar{d}_{n,p}(X, \gamma)$  is l.s.c. Furthermore, it attains a minimum and  $\bar{d}_{n,p}(X)$  is decreasing down to 0.

First we need a lemma which shows that under the strong continuity assumption made on  $\mathbb{P}_X$ , optimal (or nearly optimal), grids cannot lie in an affine hyperplane.

**Lemma 1.** Let  $p \geq 1$ . If  $\mathbb{P}_X$  is strongly continuous, then

$$\varepsilon_{d-1,p}(X) := \inf \left\{ \|\text{dist}(X, H)\|_{L^p}, H \text{ hyperplane} \right\} > 0.$$

*Proof.* Let  $\kappa > 0$  be such that  $\|\cdot\| \geq \kappa|\cdot|_2$ . Let  $H = b + u^\perp$ ,  $b \in \mathbb{R}^d$ ,  $u \in \mathbb{R}^d$ ,  $|u|_2 = 1$  (canonical Euclidean norm), be an hyperplane. If  $a \in H$ ,

$$\|X - a\| \geq \kappa|X - a|_2 \geq \kappa|(X - a, u)| = \kappa|(X - b, u)|$$

so that,  $\text{dist}(X, H) \geq \kappa|(X - b, u)|$ . Now, if  $\varepsilon_{d-1,p}(X) = 0$ , there exists two sequences  $(u_n)_{n \geq 1}$  and  $(b_n)_{n \geq 1}$  such that  $\varepsilon_n := \kappa\|(X - b_n, u_n)\|_{L^p} \rightarrow 0$ . In particular  $|(b_n, u_n)| \leq 2\|X\|_{L^p} + \varepsilon_n$ . Up to an extraction one may assume that  $u_n \rightarrow u_\infty$  (with  $|u_\infty|_2 = 1$ ) and  $(b_n, u_n) \rightarrow \ell \in \mathbb{R}$ . Then, by continuity of the  $L^p$ -norm,  $(X, u_\infty) = \ell$   $\mathbb{P}$ -a.s. which contradicts the strong continuity assumption.  $\square$

*Proof.* The proof closely follows the lines of the compactly supported case. Let  $\gamma^{(k)}$ ,  $k \geq 1$ , be a sequence of  $n$ -tuples such that  $\liminf_k \bar{d}_{n,p}(X, \gamma^{(k)}) < +\infty$ . Let  $J_\infty$  be defined like in Theorem 6 (after the appropriate extractions). Set  $\Gamma_{\gamma^{(\infty)}} = \{x_j^{(\infty)}, j \in J_\infty\}$  and  $\gamma^{(\infty)}$  accordingly.

Let  $\xi \in \mathbb{R}^d$  and let  $k'$  be a subsequence (depending on  $\xi$ ) such that  $\liminf_k \bar{F}_{n,p}(\xi, \gamma^{(k)}) = \lim_k \bar{F}_{n,p}(\xi, \gamma^{(k')})$ . We will inspect three cases:

– If  $\xi \in \limsup_k \text{Conv}\{\gamma^{(k')}\}$ , then there exists a subsequence  $k''$  such that  $\xi \in \text{Conv}\{\gamma^{(k'')}\}$  and following the lines of the proof of Theorem 6(b), one proves that either  $+\infty = \lim_k \bar{F}_{n,p}(\xi, \gamma^{(k')}) = \lim_k \bar{F}_{n,p}(\xi, \gamma^{(k'')}) \geq \bar{F}_n^p(\xi, \gamma^{(\infty)})$  or  $\xi \in \text{Conv}\{\gamma^{(\infty)}\}$  and

$$\bar{F}_{n,p}(\xi, \gamma^{(\infty)}) = F_{n,p}(\xi, \gamma^{(\infty)}) \leq \liminf_k F_{n,p}(\xi, \gamma^{(k'')}) = \lim_k \bar{F}_{n,p}(\xi, \gamma^{(k'')}) = \liminf_k \bar{F}_{n,p}(\xi, \gamma^{(k)}).$$

– If  $\xi \notin \limsup_k \text{Conv}\{\gamma^{(k')}\}$  and  $\xi \notin \partial\text{Conv}\{\gamma^{(\infty)}\}$ , then, for large enough  $k$ ,

$$\bar{F}_{n,p}(\xi, \gamma^{(k)}) = \text{dist}(\xi, \gamma^{(k)}) \rightarrow \text{dist}(\xi, \Gamma_{\gamma^{(\infty)}}) = \bar{F}_{n,p}(\xi, \gamma^{(\infty)}).$$

– Otherwise,  $\xi$  belongs to  $\partial\text{Conv}\{\gamma^{(\infty)}\}$ . At such points  $\bar{F}_{n,p}(\xi, \cdot)$  is not l.s.c. at  $\gamma^{(\infty)}$  but the boundary of the convex hull of finitely many points is made up with affine manifolds so that this boundary is negligible for  $\mathbb{P}_X$ .

Finally this proves that

$$\mathbb{P}_X(d\xi)\text{-a.s.} \quad \liminf_k \bar{F}_{n,p}(\xi, \gamma^{(k)}) \geq \bar{F}_{n,p}(\xi, \gamma^{(\infty)}).$$

One concludes using Fatou's Lemma like in the compact case that, on the one hand  $\bar{d}_{n,p}(X, \cdot)$  is l.s.c. by considering a sequence  $\gamma^{(k)}$  converging to  $\gamma^{(\infty)}$  and on the other hand that there exists an  $L^p$ -optimal grid for  $\bar{d}_{n,p}(X, \cdot)$ , namely  $\gamma^{(\infty)}$  by considering an asymptotically optimal sequence for  $(\gamma^{(k)})_{k \geq 1}$  since

$$\bar{d}_{n,p}(X) = \lim_k \bar{d}_{n,p}(X, \gamma^{(k)}) \geq \bar{d}_p(X, \Gamma_{\gamma^{(\infty)}}) \geq \bar{d}_{|\mathcal{J}_\infty|,p}(X) \geq \bar{d}_{n,p}(X)$$

so that in fact  $\bar{d}_{n,p}(X) = \bar{d}_p(X, \Gamma_{\gamma(\infty)}) = \bar{d}_{|J_\infty|,p}(X)$ .

For any grid  $\Gamma$  with size at most  $d$ ,  $\mathbb{P}(X \in \text{Conv}(\Gamma)) = 0$  so that  $\mathbb{P}_X(d\xi)$ -a.s.,  $\bar{F}_{n,p}(\xi, \Gamma) = \text{dist}(\xi, \Gamma)$  owing to the strong continuity of  $\mathbb{P}_X$ . Hence, dual and primal quantization coincide which ensures the existence of optimal grids.

Let  $n \geq d + 1$ . Assume temporarily that any optimal grids at level  $n$ , denoted  $\Gamma^{*,n}$  is “flat” i.e.  $\text{Conv}\{\Gamma^{*,n}\}$  has an empty interior or equivalently that the affine subspace spanned by  $\Gamma^{*,n}$  is included in a hyperplane  $H_n$ . Then, owing to the strong continuity assumption and Lemma 1,

$$\bar{d}_{n,p}(X) = \bar{d}^p(X, \Gamma^{*,n}) \geq \|\text{dist}(X, H_n)\|_{L^p} \geq \varepsilon_{d-1,p}(X) > 0.$$

Consequently this inequality fails for large enough  $n$  since  $\bar{d}_{n,p}(X) \rightarrow 0$  i.e.  $\overbrace{\text{Conv}\{\Gamma^{*,n}\}}^\circ$  for large enough  $n$ .

Now assume that  $\overbrace{(\text{Conv}\{\Gamma^{*,n'}\} \cap \text{supp}(\mathbb{P}_X))}^\circ \subset \Gamma^{*,n'}$  for an infinite subsequence. Let  $\xi_0$  and  $\varepsilon_0 > 0$

such that  $B(\xi_0, \varepsilon_0) \subset \text{supp}(\mathbb{P}_X)$ . This implies that  $B(\xi_0, \varepsilon_0) \cap \overbrace{\text{Conv}\{\Gamma^{*,n'}\}}^\circ = \emptyset$ . Then, for every  $\xi \in B(\xi_0, \varepsilon_0/2)$ ,  $\bar{F}_p(\xi, \Gamma^{*,n'}) = \text{dist}(\xi, \Gamma^{*,n'}) \geq (\varepsilon_0/2)$  so that

$$\bar{d}_p(X, \Gamma^{*,n'}) > (\varepsilon_0/2) \mathbb{P}(B(\xi_0, \varepsilon_0/2)) > 0$$

which contradicts the optimality of  $\Gamma^{*,n'}$  at level  $n'$  at least for  $n$  large enough. Consequently for every large enough  $n$ ,

$$\overbrace{(\text{Conv}\{\Gamma^{*,n'}\} \setminus \Gamma^{*,n'})}^\circ \cap \text{supp}(\mathbb{P}_X) \neq \emptyset.$$

Let  $\xi$  be in this nonempty set. The proof of Theorem 6(b) applies at this stage and this shows that  $\bar{d}_{n,p}(X)$  is (strictly) decreasing.  $\square$

## 5 Numerical computation of optimal dual quantizers

In order to derive optimal dual quantizers numerically, i.e. by means of gradient based optimization procedures, we have to verify the continuity of the mapping

$$\gamma \mapsto d_{n,p}(X, \gamma), \quad \gamma \in (\mathbb{R}^d)^n$$

and derive its first order derivative.

Therefore, we will need the assumption of dual non-degeneracy in the Linear Program  $F_n^p(\xi, \gamma)$  to establish the gradient of  $d_n^p(X, \cdot)$ . We therefore call a grid  $\Gamma_\gamma = \{x_1, \dots, x_n\}$  non-degenerated, if for every  $I \in \mathcal{I}(\Gamma_\gamma)$  and  $\mathbb{P}_X$ -a.e.  $\xi \in D_I \cap \text{supp}(\mathbb{P}_X)$  it holds

$$A_{I^c}^T u < c_{I^c}, \quad \text{where } u = (A_I^T)^{-1} c_I.$$

This condition implies together with the primal non-degeneracy  $\lambda_I = A_I^{-1} b > 0$  the uniqueness of the primal and dual solutions for (LP).

In the Euclidean case e.g., this assumption is fulfilled as soon as the Delaunay triangulation is non-degenerated, i.e. no  $d+2$  points lie on a hypersphere, which then also implies the uniqueness of the Delaunay triangulation.

**Theorem 8.** *Let  $X \in L^p(\mathbb{P})$ ,  $p \geq 1$  and assume that  $\mathbb{P}_X$  satisfies the strong continuity assumption. Moreover, let  $\gamma_0 = (x_1, \dots, x_n)$  be a  $n$ -tuple in  $(\mathbb{R}^d)^n$  such that  $\text{supp}(\mathbb{P}_X) \subset \text{conv}\{\Gamma_{\gamma_0}\}$ . Then*

(a) the mapping

$$\gamma \mapsto d_{n,p}(X, \gamma), \quad \gamma \in (\mathbb{R}^d)^n$$

is continuous in  $\gamma_0$ .

(b) If furthermore  $x \mapsto \|x\|^p$  is differentiable and  $\gamma_0 = (x_1, \dots, x_n)$  is non-degenerated in the above sense, then  $d_n^p(X, \cdot)$  is differentiable at  $\gamma_0$  with derivative

$$\frac{\partial}{\partial x_i^j} d_n^p(X, \gamma_0) = \mathbb{E} \left[ \lambda_i(X) \left( \frac{\partial}{\partial j} \|X - x_i\|^p - u_j(X) \right) \right], \quad 1 \leq j \leq d, 1 \leq i \leq n,$$

where  $\lambda(X)$  and  $u(X)$  are the  $\mathbb{P}_X$ -a.s. unique primal and dual solutions for the Linear Program  $F_n^p(X, \gamma_0)$ .

*Proof.* (a) Due to Theorem 6(a), it remains to show that  $d_n^p(X, \cdot)$  is u.s.c. at  $\gamma_0 = (x_1, \dots, x_n)$ . Therefore, denote by  $H_{\gamma_0}$  the set of all hyperplanes generated by any subset  $\{x_{i_1}, \dots, x_{i_d}\}$  of  $\Gamma_{\gamma_0}$  and let  $\gamma_k = (x_1^k, \dots, x_n^k) \in (\mathbb{R}^d)^n$  be a sequence converging to  $\gamma_0$  for  $k \rightarrow \infty$ . We will then show for every  $\xi \in \text{supp}(\mathbb{P}_X) \setminus H_{\gamma_0}$

$$\limsup_{k \rightarrow \infty} F_n^p(X, \gamma_k) \leq F_n^p(\xi, \gamma_0).$$

Consequently, let  $\xi \in \text{supp}(\mathbb{P}_X) \setminus H_{\gamma_0}$  and choose a basis  $I \in \mathcal{I}(\Gamma_{\gamma_0})$  such that  $\xi \in D_I(\Gamma_{\gamma_0})$ . Since  $\xi \notin H_{\gamma_0}$ , it lies in the interior of  $\text{conv}\{x_j : j \in I\}$ , which implies  $\lambda_I = A_I^{-1}b > 0$  and

$$F_n^p(\xi, \gamma_0) = \lambda_I^T c_I.$$

Denoting

$$A^k = \begin{bmatrix} x_1^k & \dots & x_n^k \\ 1 & \dots & 1 \end{bmatrix}, \quad c = \begin{bmatrix} \|\xi - x_1^k\|^p \\ \vdots \\ \|\xi - x_n^k\|^p \end{bmatrix},$$

we clearly have  $A^k \rightarrow A$  and  $c^k \rightarrow c$  as  $k \rightarrow \infty$ .

Moreover,  $A_I^k$  is regular for  $k$  large enough, so that it also holds  $(A_I^k)^{-1} \rightarrow A_I^{-1}$ . But this also implies for  $\lambda^k = (A_I^k)^{-1}b$

$$\lambda^k \rightarrow \lambda_I \quad \text{and} \quad \lambda > 0 \quad \text{for } k \text{ large enough.}$$

Therefore,  $\lambda^k$  becomes a feasible solution for  $F_n^p(\xi, \gamma_k)$ , which yields

$$\limsup_{k \rightarrow \infty} F_n^p(\xi, \gamma_k) \leq \lim_{k \rightarrow \infty} (\lambda^k)^T c^k = \lambda_I^T c_I = F_n^p(\xi, \gamma_0).$$

Since  $\mathbb{P}(X \in H_{\gamma_0}) = 0$  and  $d_n^p(X, \gamma_0) < +\infty$  by assumption, Fatou's Lemma yield the u.s.c. of  $d_n^p(X, \cdot)$  in  $\gamma_0$ .

(b) Denote by  $N_{\gamma_0}$  the set where  $F_n^p(\xi, \gamma_0)$  is degenerated in the dual sense, which is by assumption a null set and moreover let  $\xi \in \text{supp}(\mathbb{P}_X) \setminus (H_{\gamma_0} \cup N_{\gamma_0})$ . Then the Linear Program  $F_n^p(\xi, \gamma_0)$  is non-degenerated in the primal and dual sense, so that it is classical background from Linear Programming theory, that there is a unique  $I \in \mathcal{I}(\Gamma_{\gamma_0})$  such that  $\lambda_I = A_I^{-1}b$  and  $u = (A_I^T)^{-1}c_I$  are the unique solutions for  $F_n^p(\xi, \gamma_0)$ , i.e.

$$F_n^p(\xi, \gamma_0) = \lambda_I^T c_I = u^T b. \quad (19)$$

Moreover, one checks under these assumptions that after reordering of rows for  $\lambda = (\lambda_i, 0)$  it holds

$$c - A^T u + \lambda > 0.$$



Since

$$\gamma \mapsto c - A^T u + \lambda$$

is continuous in the point  $\gamma_0$ , there exists a neighborhood  $\mathcal{U}(\gamma_0)$  of  $\gamma_0$  such that for every  $\bar{\gamma} = (\bar{x}_1, \dots, \bar{x}_n) \in \mathcal{U}(\gamma_0)$

$$\bar{c} - \bar{A}^T \bar{u} + \bar{\lambda} > 0.$$

for

$$\bar{A} = \begin{bmatrix} \bar{x}_1 & \dots & \bar{x}_n \\ 1 & \dots & 1 \end{bmatrix}, \quad c = \begin{bmatrix} \|\xi - \bar{x}_1\|^p \\ \vdots \\ \|\xi - \bar{x}_n\|^p \end{bmatrix}, \quad \bar{\lambda} = (\bar{A}_I^{-1} b, 0), \quad \bar{u} = (\bar{A}_I^T)^{-1} \bar{c}_I.$$

But this implies by Proposition 5 that  $I$  is also optimal for every  $\bar{\gamma} \in \mathcal{U}(\gamma_0)$  so that we conclude

$$F_n^p(\xi, \bar{\gamma}) = \bar{\lambda}_I^T \bar{c}_I = \bar{u}^T b.$$

Therefore we may differentiate the identity (19) formally with respect to  $\gamma_0$  and obtain

$$\nabla_{\gamma_0} (F_n^p(\xi, \gamma_0)) = \nabla(c_I^T \lambda_i) = \nabla(c_I^T A_I^{-1} b) = (\nabla c_I)^T A_I^{-1} b + c_I^T \nabla(A_I^{-1}) b.$$

Using the identity  $\nabla(A^{-1}) = -A^{-1} \nabla A A^{-1}$  we derive

$$\begin{aligned} \nabla_{\gamma_0} (F_n^p(\xi, \gamma_0)) &= (\nabla c_I)^T \lambda_I - c_I^T A_I^{-1} \nabla A_I A_I^{-1} b \\ &= (\nabla c_I)^T \lambda_I - u^T \nabla A_I \lambda_I. \end{aligned}$$

Some elementary tensor calculus then yields

$$\frac{\partial}{\partial x_i^j} F_n^p(\xi, \gamma_0) = \lambda_i(\xi) \left( \frac{\partial}{\partial j} \|\xi - x_i\|^p - u_i^j(\xi) \right),$$

which is bounded as a function of  $\xi$  on any compact set, so that the assertion follows.  $\square$

## 5.1 One dimensional setting

In the one dimensional case, we can derive, due to a simpler geometrical structure, more explicit expressions for  $F_n^p$  and its derivatives.

To be more precisely, let  $\gamma = (x_1, \dots, x_n) \in \mathbb{R}$  be ordered non-decreasingly. Then

$$D_I(\Gamma_\gamma) = [x_i, x_{i+1}] \quad \text{for } I = \{i, i+1\},$$

so that we arrive at a dual quantization error

$$d_n^p(X, \gamma) = \sum_{i=1}^{n-1} \frac{1}{x_{i+1} - x_i} \int_{x_i}^{x_{i+1}} (x_{i+1} - \xi)(\xi - x_i)^p + (\xi - x_i)(x_{i+1} - \xi)^p \mathbb{P}_X(d\xi). \quad (20)$$

**Uniform distribution:** For the uniform distribution  $\mathcal{U}([0, 1])$  we can even compute the exact solutions for the dual quantization problem. Therefore, one easily derives from (20)

$$d_n^p(\mathcal{U}([0, 1]), \gamma) = \frac{2}{(p+1)(p+2)} \sum_{i=1}^{n-1} (x_{i+1} - x_i)^{p+1},$$

so that setting  $y_i = x_{i+1} - x_i$  yields

$$d_n^p(\mathcal{U}([0, 1])) = \frac{2}{(p+1)(p+2)} \min \left\{ \sum_{i=1}^{n-1} y_i^{p+1} : \sum_i y_i = 1, y_i \geq 0 \right\},$$

where we have to fix the grid endpoints  $x_1 = 0$  and  $x_n = 1$  to ensure  $[0, 1] \subset \text{conv}\{\Gamma_\gamma\}$ . The solution to this problem is obviously given by  $y_i = \frac{1}{n-1}$ , which implies

$$x_i^* = \frac{i-1}{n-1} \quad \text{and} \quad d_n^p(\mathcal{U}([0, 1])) = \frac{2}{(p+1)(p+2)} \frac{1}{(n-1)^p}.$$

Recall that it holds for ordinary quantization of the uniform distribution

$$x_i^{*,\text{vq}} = \frac{2i-1}{2n} \quad \text{and} \quad e_n^p(\mathcal{U}([0, 1])) = \frac{1}{2^p(p+1)} \frac{1}{n^p},$$

so that we conclude for the sharp asymptotics

$$\lim_{n \rightarrow \infty} n^{1/d} d_n^p(\mathcal{U}([0, 1])) = \left( \frac{2^{p+1}}{p+2} \right)^{1/p} \lim_{n \rightarrow \infty} n^{1/d} e_n^p(\mathcal{U}([0, 1])).$$

Furthermore, we recognize that an optimal dual quantizer of size  $n+1$  is made up by the midpoints of an optimal regular quantizer of size  $n$  plus the interval endpoints. One may even show in this context that such a construction leads to asymptotically optimal dual quantizers for any compactly supported distribution in dimension one.

**General quadratic case:** In the general quadratic setup, we derive from Theorem 8 for  $p = 2$  and an ordered grid  $\gamma = (x_1, \dots, x_n)$

$$\frac{\partial}{\partial x_i} d_n^p(X, \gamma) = \int_{x_{i-1}}^{x_{i+1}} \xi \mathbb{P}_X(d\xi) - x_{i-1} \int_{x_{i-1}}^{x_i} \mathbb{P}_X(d\xi) - x_{i+1} \int_{x_i}^{x_{i+1}} \mathbb{P}_X(d\xi), \quad 2 \leq i \leq n-1.$$

If  $\text{conv}\{\text{supp}(\mathbb{P}_X)\} = [a, b]$ , we statically fix the endpoints  $x_1 = a$  and  $x_n = b$  in any optimization procedure to generate optimal dual quantizers. Otherwise, in the unbounded case, we introduce boundary conditions according to a nearest neighbor mapping

$$\begin{aligned} \frac{\partial}{\partial x_1} \bar{d}_n^p(X, \gamma) &= 2 \int_{-\infty}^{x_1} (x_1 - \xi) \mathbb{P}_X(d\xi) + \int_{x_1}^{x_2} (\xi - x_2) \mathbb{P}_X(d\xi) \\ \frac{\partial}{\partial x_n} \bar{d}_n^p(X, \gamma) &= 2 \int_{x_n}^{+\infty} (x_n - \xi) \mathbb{P}_X(d\xi) + \int_{x_{n-1}}^{x_n} (\xi - x_{n-1}) \mathbb{P}_X(d\xi). \end{aligned}$$

The second derivative then reads for absolutely continuous  $\mathbb{P}_X$

$$\begin{aligned} \frac{\partial^2}{\partial^2 x_1} \bar{d}_n^p(X, \gamma) &= 2 \int_{-\infty}^{x_1} \mathbb{P}_X(d\xi) - (x_2 - x_1) \frac{d\mathbb{P}_X}{d\lambda^1}(x_1) \\ \frac{\partial^2}{\partial x_2 \partial x_1} \bar{d}_n^p(X, \gamma) &= \frac{\partial^2}{\partial x_1 \partial x_2} d_n^p(X, \gamma) = - \int_{x_1}^{x_2} \mathbb{P}_X(d\xi) \\ \frac{\partial^2}{\partial^2 x_i} d_n^p(X, \gamma) &= (x_{i+1} - x_{i-1}) \frac{d\mathbb{P}_X}{d\lambda^1}(x_i), \quad 2 \leq i \leq n-1 \\ \frac{\partial^2}{\partial x_{i+1} \partial x_i} d_n^p(X, \gamma) &= \frac{\partial^2}{\partial x_i \partial x_{i+1}} d_n^p(X, \gamma) = - \int_{x_i}^{x_{i+1}} \mathbb{P}_X(d\xi), \quad 2 \leq i \leq n-1 \\ \frac{\partial^2}{\partial x_{n-1} \partial x_n} \bar{d}_n^p(X, \gamma) &= \frac{\partial^2}{\partial x_n \partial x_{n-1}} d_n^p(X, \gamma) = - \int_{x_{n-1}}^{x_n} \mathbb{P}_X(d\xi) \\ \frac{\partial^2}{\partial^2 x_n} \bar{d}_n^p(X, \gamma) &= 2 \int_{x_n}^{+\infty} \mathbb{P}_X(d\xi) - (x_n - x_{n-1}) \frac{d\mathbb{P}_X}{d\lambda^1}(x_n). \end{aligned}$$

The above integral expressions can be for most distributions evaluated in closed-form. Therefore, it is straightforward to employ a Newton method to find a zero of  $\nabla d_n^p(X, \cdot)$ , which yields an optimal dual quantizer. Such a procedure, initialized with an equidistant grid in the center of the distribution, converges usually very fast ( $< 10$  iterations) to an optimal grid.

## 5.2 Multi-dimensional setting

In the multi-dimensional case, the computation of  $\nabla d_n^p(X, \cdot)$  involves the evaluation of multi-dimensional integrals, for which in general no closed-form solution is available and numerical evaluation of these integrals is a rather time consuming task.

We therefore focus, as in the case of regular quantization, on a Robbins-Monro stochastic optimization algorithm. Such an algorithm has the advantage of building up the necessary gradient information step-by-step during the simulation and therefore is by several magnitudes faster than a “batch”-approach which evaluates the full gradient at each iteration.

This variant of the Robbins-Monro algorithm is in the case of regular vector quantization also known as *Competitive Vector Learning Quantization* algorithm (CVLQ) (see [10]).

---

### Algorithm 1 CVLQ for dual Quantization

---

**Input:**

- Step sequence  $\alpha_k \geq 0$  such that  $\sum_{k \geq 0} \alpha_k = +\infty$ ,  $\sum_{k \geq 0} \alpha_k^2 < +\infty$
- Initial grid  $\gamma_0 \in (\mathbb{R}^d)^n$

**Main loop:**

```

for  $k = 0$  to  $N - 1$  do
  Generate i.i.d. sample  $X_k \sim X$ 
  Set
     $\gamma_{k+1} \leftarrow \gamma_k - \alpha_k \nabla_{\gamma_k} F_n^p(X_k, \gamma_k)$ 

```

**end for**

---

To compare this procedure to the regular CVLQ-algorithm, we inspect the main loop for the case  $p = 2$ . Given a realization  $X_k$  of  $X$ , we only have to replace the Nearest Neighbor search by a search for the Delaunay triangle  $I^*$ , which contains  $X_k$ . According to Theorem 5, the primal solution  $\lambda_I^*$  to the Linear Program  $F_n^p(X_k, \gamma)$  is then given by the barycentric coordinates of  $X_k$  in the triangle  $I^*$  and the dual solution can be calculated by the formula

$$u^* = 2(z^* - X_k),$$

where  $z^*$  is the center of the hypersphere spanning the triangle  $I^*$ . We therefore can simplify the partial derivative of  $F_n^p(X_k, (x_1, \dots, x_n))$  for  $I^*$  being the Delaunay triangle containing  $X_k$  to

$$\frac{\partial}{\partial x_i} F_n^p(X_k, (x_1, \dots, x_n)) = 2\lambda_i^*(x_i - z^*).$$

Main loop:: regular CVLQ	Main loop:: CVLQ for dual quantization
<b>for</b> $k = 0$ to $N - 1$ <b>do</b>	<b>for</b> $k = 0$ to $N - 1$ <b>do</b>
Generate i.i.d. sample $X_k \sim X$	Generate i.i.d. sample $X_k \sim X$
Find NN index $i^*$ of $X_k$ in $(x_1^k, \dots, x_n^k)$	Find Delaunay triangle $I^*$ in $(x_1^k, \dots, x_n^k)$ , which contains $X_k$
	Compute LP solution $\lambda_I^*$ and center $z^*$
<b>for</b> $j = 1$ to $n$ <b>do</b>	<b>for</b> $j = 1$ to $n$ <b>do</b>
<b>if</b> $j = i^*$ <b>then</b>	<b>if</b> $j \in I^*$ <b>then</b>
$x_j^{k+1} \leftarrow x_j^k - \alpha_k \cdot (x_j^k - X_k)$	$x_j^{k+1} \leftarrow x_j^k - \alpha_k \cdot \lambda_j^* \cdot (x_j^k - z^*)$
<b>else</b>	<b>else</b>
$x_j^{k+1} \leftarrow x_j^k$	$x_j^{k+1} \leftarrow x_j^k$
<b>end if</b>	<b>end if</b>
<b>end for</b>	<b>end for</b>
<b>end for</b>	<b>end for</b>

These procedures usually converge quickly to a first approximation of an optimal quantization grid. For a local refinement, we propose to combine the above approach with a few quasi-Newton steps of a deterministic optimization algorithm, where the evaluation of the integral expression is performed by a Monte Carlo- resp. Quasi Monte Carlo method (cf. [16]).

Numerical results of this approach are given for the Uniform distribution on  $[0, 1]^2$  in figures 1 to 4 with grid sizes 8 to 16 and for the standard normal distribution on  $\mathbb{R}^2$  for a grid size of 250 in figure 5.

## References

- [1] V. Bally and G. Pagès. A quantization algorithm for solving multi-dimensional discrete-time optimal stopping problems. *Bernoulli*, 9(6):1003–1049, 2003.
- [2] V. Bally, G. Pagès and J. Printems. A quantization tree method for pricing and hedging multidimensional american options. *Mathematical Finance*, 15:119–168(50), January 2005.
- [3] O. Bardou, S. Bouthemy and G Pagès. Optimal Quantization for the Pricing of Swing Options. *Applied Mathematical Finance*, 16(2):183–217, 2009.
- [4] O. Bardou, S. Bouthemy and G Pagès. When are Swing options bang-bang? *International Journal of Theoretical and Applied Finance (IJTAF)*, 13(06):867–899, 2010.
- [5] A. L. Bronstein, G. Pagès and B. Wilbertz. How to speed up the quantization tree algorithm with an application to swing options. *Quantitative Finance*, 10(9):995 – 1007, November 2010.
- [6] E. Gobet, G. Pagès, H. Pham and J. Printems. Discretization and simulation of the Zakai equation. *SIAM J. Numer. Anal.*, 44(6):2505–2538 (electronic), 2006.
- [7] S. Graf and H. Luschgy. *Foundations of Quantization for Probability Distributions*. Lecture Notes in Mathematics n<sup>o</sup>1730. Springer, Berlin, 2000.
- [8] M. Padberg. *Linear optimization and extensions*, volume 12 of *Algorithms and Combinatorics*. Springer-Verlag, Berlin, 1995.

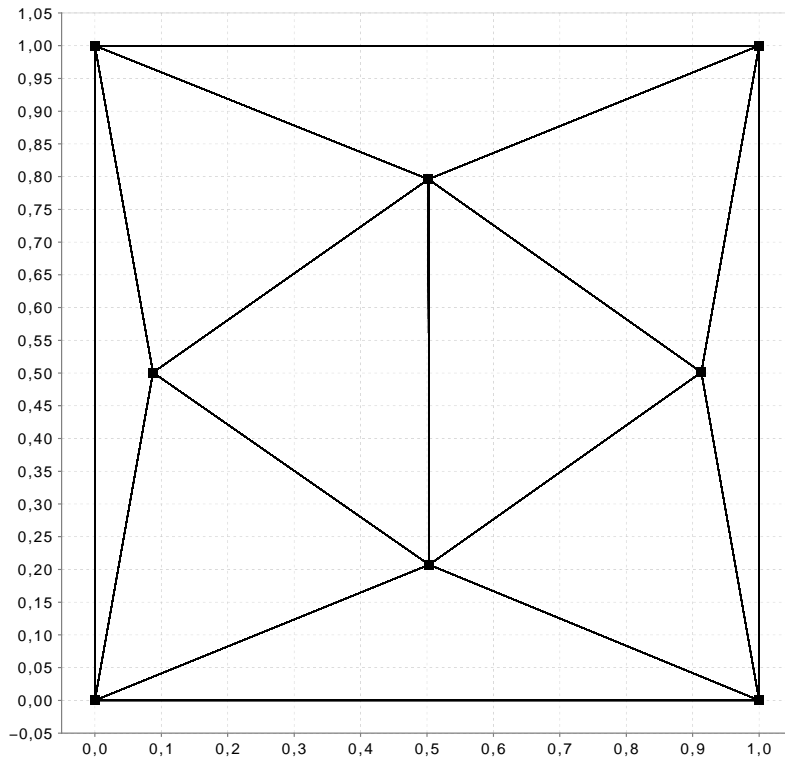


Figure 1: Dual Quantization for  $\mathcal{U}([0, 1]^2)$  and  $N = 8$

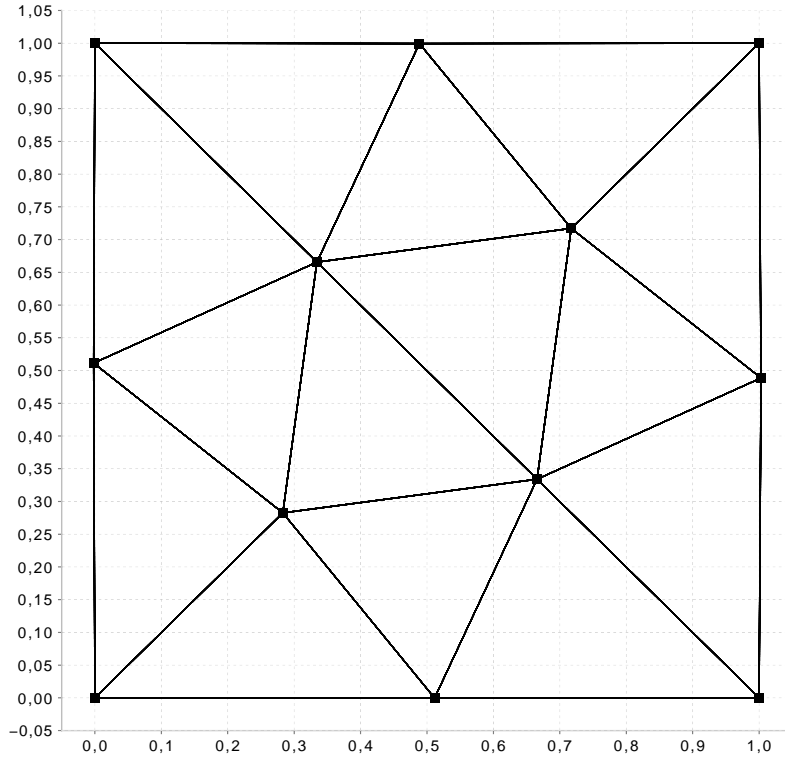


Figure 2: Dual Quantization for  $\mathcal{U}([0, 1]^2)$  and  $N = 12$

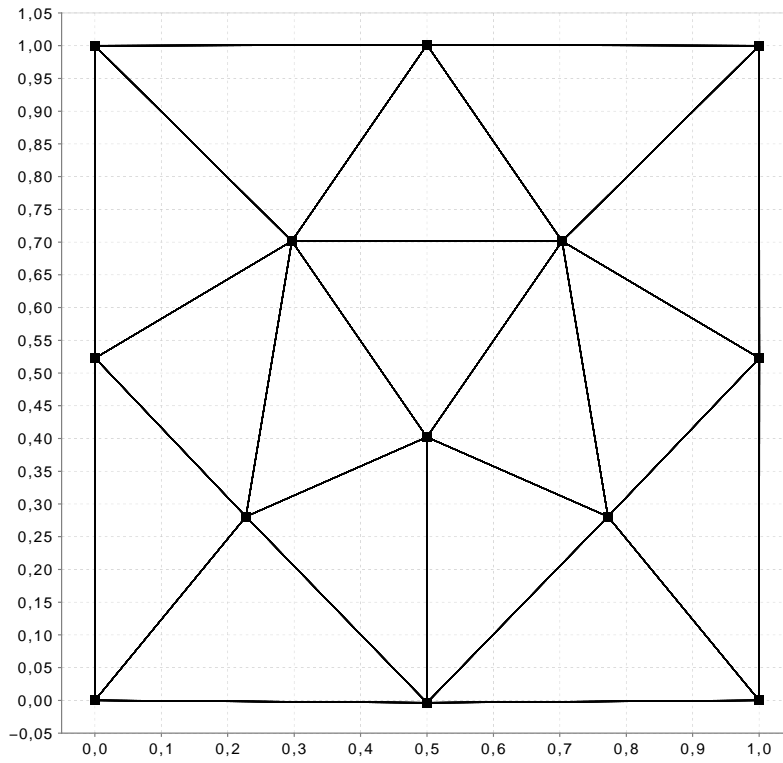


Figure 3: Dual Quantization for  $\mathcal{U}([0, 1]^2)$  and  $N = 13$

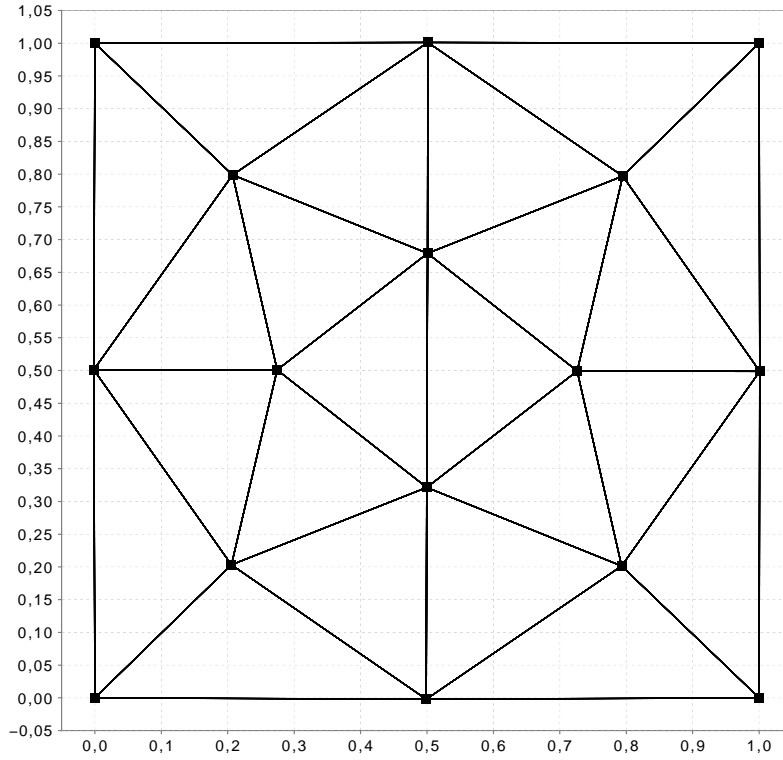


Figure 4: Dual Quantization for  $\mathcal{U}([0, 1]^2)$  and  $N = 16$

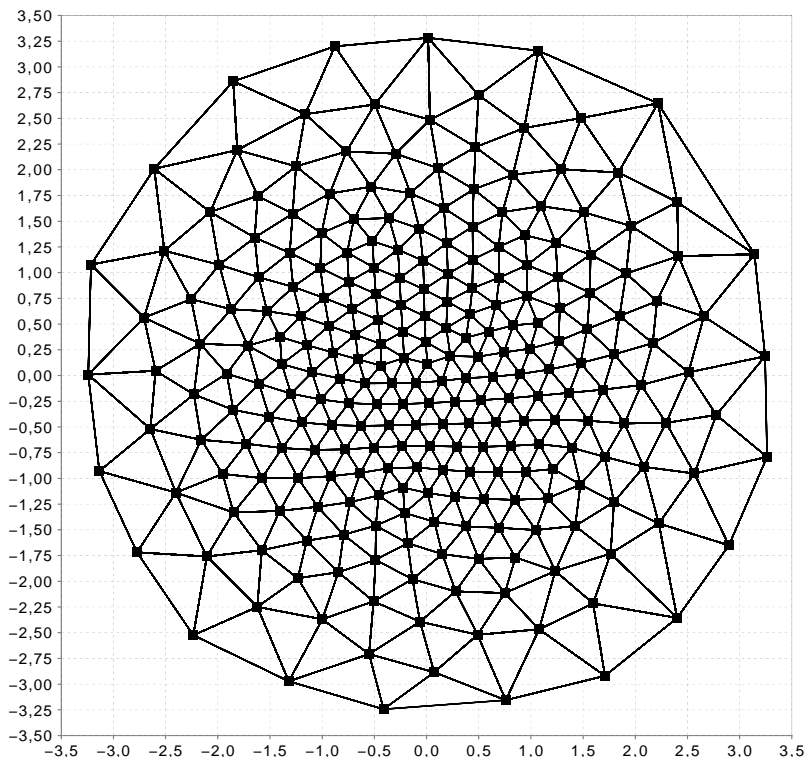


Figure 5: Dual Quantization for  $\mathcal{N}(0, I_2)$  and  $N = 250$

- [9] G. Pagès, H. Pham and J. Printems. Optimal quantization methods and applications to numerical methods and applications in finance. In S. Rachev, editor, *Handbook of Computational and Numerical Methods in Finance*, pages 253–298. Birkhäuser, 2004.
- [10] G. Pagès and J. Printems. Optimal quadratic quantization for numerics: the gaussian case. *Monte Carlo Methods and Applications*, 9(2):135–166, 2003.
- [11] G. Pagès and B. Wilbertz. Sharp rate for the dual quantization problem. In progress, 2010.
- [12] H. Pham, W. Runggaldier and A. Sellami. Approximation by quantization of the filter process and applications to optimal stopping problems under partial observation. *Monte Carlo Methods Appl.*, 11(1):57–81, 2005.
- [13] V. T. Rajan. Optimality of the delaunay triangulation in  $R^d$ . In *SCG '91: Proceedings of the seventh annual symposium on Computational geometry*, pages 357–363, New York, NY, USA, 1991. ACM.
- [14] A. Sellami. Comparative survey on nonlinear filtering methods: the quantization and the particle filtering approaches. *J. Stat. Comput. Simul.*, 78(1-2):93–113, 2008.
- [15] A. Sellami. Quantization based filtering method using first order approximation. *SIAM J. Numer. Anal.*, 47(6):4711–4734, 2010.
- [16] B. Wilbertz. Computational aspects of Functional Quantization for Gaussian measures and applications. Diploma Thesis, Trier University, 2005.