

# An approach to Stochastic Process using Quasi-Arithmetic Means

Etienne Cuvelier and Monique Noirhomme-Fraiture

Institut d'Informatique (FUNDP)  
21, rue Grandgagnage, 5000 Namur, Belgium  
(e-mail: [ecu@info.fundp.ac.be](mailto:ecu@info.fundp.ac.be), [mno@info.fundp.ac.be](mailto:mno@info.fundp.ac.be))

**Abstract.** Probability distributions are central tools for probabilistic modeling in data mining. In functional data analysis (FDA) they are weakly studied in the general case. In this paper we discuss a probability distribution law for functional data considered as stochastic process. We define first a new kind of stationarity linked to the Archimedean copulas, and then we build a probability distribution using jointly the Quasi-arithmetic means and the generators of Archimedean copulas. We also study some properties of this new mathematical tool.

**Keywords:** Functional Data Analysis, Probability distributions, Stochastic Process, Quasi-Arithmetic Mean, Archimedean copulas.

## 1 Introduction

Probability distributions are central tools for probabilistic modeling in data mining. In functional data analysis, as functional random variable can be considered as stochastic process, the probability distribution have been studied largely, but with rather strong hypotheses, [Cox and Miller, 1965], [Gihman and Skorohod, 1974], [Bartlett, 1978] and [Stirzaker, 2005]. Some processes are very famous like Markov process [Meyn and L, 1993]. Such a process has the property that present is not influenced by all the past but only by the last visited state. A very particular case is the random walk, which has the property that one-step transitions are permitted only to the nearest neighboring states. Such local changes of state may be regarded as the analogue for discrete states of the phenomenon of continuous changes for continuous states. The limiting process is called the Wiener process or Brownian motion. The Wiener process is a diffusion process having the special property of independent increments. Some more general Markov chain with only local changes of state are permissible, gives also Markov limiting process for continuous time and continuous states. The density probability is solution of a special case of the Fokker-Planck diffusion equation.

In preceding work [Cuvelier and Noirhomme-Fraiture, 2005] we used copulas to model the distribution of functional random variables at discrete cutting points. Here, using the separability concept, we can consider the continuous case as the limit of the discrete one. We will use quasi-arithmetic means in order to avoid copulas problem when considering the limit when the number

of cuttings tends to infinity.

In section 2 we define the concept of distribution of functions and recall the notion of separability. In section 3 we propose to use the Quasi-arithmetic mean in conjunction with an Archimedean generator to build a probability distributions appropriate to the dimensional infinite nature of the functional data. And in section 4 we study the properties of this new mathematical tool.

## 2 Distribution of a functional random variable

Let us recall some definitions that will be useful in the following paper.

**Definition 1.** Let  $(\Omega, \mathcal{A}, P)$  a probability space and  $\mathcal{D}$  a closed real interval. A *functional random variable (frv)* is any function from  $\mathcal{D} \times \Omega \rightarrow \mathbb{R}$  such for any  $t \in \mathcal{D}$ ,  $X(t, \cdot)$  is a real random variable on  $(\Omega, \mathcal{A}, P)$ . Each function  $X(\cdot, \omega)$  is called a realization. In the following we will write  $\underline{X}$  for  $X(\cdot, \omega)$ , and  $\underline{X}_t$  for  $X(t, \cdot)$ .  $\underline{X}_t$  can be considered as a stochastic process.

We study, here, the measurable and bounded functions.

**Definition 2.** Let  $\mathcal{D}$  a closed real interval, then  $\mathcal{L}_2(D)$  is the space of real measurable functions  $u(t)$  defined on a real interval  $\mathcal{D}$  such that

$$\|u\|_2 = \left\{ \int_{\mathcal{D}} |u(t)|^2 dt \right\}^{1/2} < \infty \quad (1)$$

**Definition 3.** Let  $f, g \in \mathcal{L}_2(D)$ . The pointwise order between  $f$  and  $g$  on  $\mathcal{D}$  is defined as follows :

$$\forall t \in \mathcal{D}, f(t) \leq g(t) \iff f \leq_{\mathcal{D}} g \quad (2)$$

**Definition 4.** The *functional cumulative distribution function (fcdf)* of a frv  $\underline{X}$  on  $\mathcal{L}_2(D)$  computed at  $u \in \mathcal{L}_2(D)$  is given by :

$$F_{\underline{X}, \mathcal{D}}(u) = P\{\underline{X} \leq_{\mathcal{D}} u\} \quad (3)$$

**Definition 5.** A frv is called *separable* if there exists in  $\mathcal{D}$  an everywhere countable set  $I$  of points  $\{t_i\}$  and a set  $N$  of  $\Omega$  of probability 0 such that for an arbitrary open set  $G \subset \mathcal{D}$  and an arbitrary closed set  $F \subset \mathbb{R}$  the two sets

$$\begin{aligned} &\{\omega : X(t, \omega) \in F, \forall t \in G\} \\ &\{\omega : X(t, \omega) \in F, \forall t \in G \cap I\} \end{aligned}$$

differ from each other only on the subset  $N$ . The set  $I$  is called the separability set [Gihman and Skorohod, 1974].

The space  $\mathcal{L}_2(D)$  is a separable Hilbert space. In the following we suppose that any realization of  $\underline{X}$  is in  $\mathcal{L}_2(\mathcal{D})$ .

**Definition 6.** Two frv  $X_1(t, \omega)$  and  $X_2(t, \omega)$  ( $t \in \mathcal{D}, \omega \in \Omega$ ) are called stochastically equivalent if for any  $t \in \mathcal{D}$

$$P \{X_1(t, \omega) \neq X_2(t, \omega)\} = 0 \tag{4}$$

The interest of separability comes from the following theorem .

**Theorem 1 (J.L. Doob).** Let  $\mathcal{X}$  and  $\mathcal{Y}$  be metric spaces,  $\mathcal{X}$  be separable,  $\mathcal{Y}$  be compact. An arbitrary random function  $X(t, \omega)$ ,  $t \in \mathcal{X}$  with values in  $\mathcal{Y}$  is stochastically equivalent to a certain separable random function.

### 3 The QAMM and QAMML distributions

In this section we build a sequence of sets that converge toward a separability set of  $\mathcal{D}$  and at each step we define a probability distribution. Let  $n \in \mathbb{N}$ , and  $\{t_1^n, \dots, t_n^n\}$ ,  $n$  equidistant points of  $\mathcal{D}$  such that  $t_1^n = \inf(\mathcal{D})$  and  $t_n^n = \sup(\mathcal{D})$ , and  $\forall i \in \{1, \dots, n - 1\}$  we have  $|t_{i+1}^n - t_i^n| = \frac{|\mathcal{D}|}{n} = \Delta_t$ . Let the two following sets

$$A_n(u) = \bigcap_{i=1}^n \{\omega \in \Omega : X(t_i^n, \omega) \leq u(t_i^n)\}$$

$$\mathcal{A}(u) = \{\omega \in \Omega : \underline{X} \leq_{\mathcal{D}} u\}$$

We will use the following distribution to approximate the *fcdf* (3):

$$P[A_n(u)] = H(u(t_1^n), \dots, u(t_n^n)) \tag{5}$$

where  $H(\cdot, \dots, \cdot)$  is a joint distribution of dimension  $n$ . In previous works (see [Diday, 2002], [Vrac *et al.*, 2001], [Cuvelier and Noirhomme-Fraiture, 2005]) the Archimedean copulas were used for the approximation with small value of  $n$ . Let us recall the definition and property of copulas.

**Definition 7.** A copula is a multivariate cumulative distribution function defined on the  $n$ -dimensional unit cube  $[0, 1]^n$  such that every marginal distribution is uniform on the interval  $[0, 1]$  :

$$C : [0, 1]^n \rightarrow [0, 1] : (u_1, \dots, u_n) \mapsto C(u_1, \dots, u_n)$$

The power of copulas comes from the following theorem (see [Nelsen, 1999]).

**Theorem 2 (Sklar’s theorem).** Let  $H$  be an  $n$ -dimensional distribution function with margins  $F_1, \dots, F_n$ . Then there exists an  $n$ -copula  $C$  such that for all  $x \in \mathbb{R}^n$ ,

$$H(x_1, \dots, x_n) = C(F_1(x_1), \dots, F_n(x_n)). \tag{6}$$

If  $F_1, \dots, F_n$  are all continuous, then  $C$  is unique; otherwise,  $C$  is uniquely determined on Range of  $F_1 \times \dots \times$  Range of  $F_n$ .

Before using copulas, we define a function that gives the distribution of the values of  $\underline{X}_t$  for a chosen  $t \in \mathcal{D}$ .

**Definition 8.** Let  $\underline{X}$  a frv. We define the *surface of distributions* as follow :

$$G(t, y) = P\{\underline{X}_t \leq y\} \tag{7}$$

We can use various methods for determining suitable  $G$  for a chosen value of  $t$ . Thus for example, if  $\underline{X}$  is a Gaussian process with mean value  $\mu(t)$  and standard deviation  $\sigma(t)$ , then we can use the *cdf* from  $\mathcal{N}(\mu(t), \sigma(t))$ . In other cases we can use the empirical cumulative distribution function to estimate  $\hat{G}$  :

$$\hat{G}(t, y) = \frac{\#\{X_i(t) \leq y\}}{N} \tag{8}$$

In the following we will always use this function  $G$  in conjunction with a function  $u$  of  $\mathcal{L}_2(\mathcal{D}) : G[t, u(t)]$ . So, for ease the notations, we will write  $G[t; u] = G[t, u(t)]$ . If we use the preceding expression in conjunction with (6), then (5) become :

$$P[\mathcal{A}_n(u)] = C(G[t_1^n; u], \dots, G[t_q^n; u]) \tag{9}$$

An important class of stochastic process is the class of stationary processes. A stochastic process is said to be *strictly stationary* [Burril, 1972] if its distributions do not change with time; i.e. if for any  $t_1, \dots, t_n \in \mathcal{D}$  and for any  $h \in \mathcal{D}$ , the multivariate distribution function of  $(\underline{X}_{t_1+h}, \dots, \underline{X}_{t_n+h})$  does not depend on  $h$ . We propose here a more wide stationary property.

**Definition 9.** A stochastic process is said *copula stationary* if  $\forall t_1, \dots, t_n \in \mathcal{D}$  and for any  $h \in \mathcal{D}$ , the copula of  $(\underline{X}_{t_1+h}, \dots, \underline{X}_{t_n+h})$  does not depend on  $h$ , i.e. its copula does not change with time.

Let us notice that, if we deal with true functional data, realizations of a stochastic process  $\underline{X}$ , we can suppose that there is always the same functional relation between  $\underline{X}_s$  and  $\underline{X}_t$  for any value  $s, t \in \mathcal{D}$ . If a frv is also a *copula stationary* stochastic process, then we call it a *copula stationary frv*. There is an important class of copulas which is well appropriate for *copula stationary* stochastic processes : the class of Archimedean copulas.

**Definition 10.** An Archimedean copula is a function from  $[0, 1]^n$  to  $[0, 1]$  given by

$$C(u_1, \dots, u_n) = \psi \left[ \sum_{i=1}^n \phi(u_i) \right] \tag{10}$$

where  $\phi$ , called the generator, is a function from  $[0, 1]$  to  $[0, \infty]$  such that:

- $\phi$  is a continuous strictly decreasing function,
- $\phi(0) = \infty$  and  $\phi(1) = 0$ ,

**Table 1.** Families of completely monotonic generators

Name	Generator	Dom. of $\theta$
Clayton	$t^\theta - 1$	$\theta > 0$
Frank	$-\ln \frac{e^{-\theta t} - 1}{e^{-\theta} - 1}$	$\theta > 0$
Gumbel-Hougaard	$(-\ln t)^\theta$	$\theta \geq 1$

- $\psi = \phi^{-1}$  is completely monotonic on  $[0, \infty[$  i.e.  $(-1)^k \frac{d^k}{dt^k} \psi(t) \geq 0$  for all  $t$  in  $[0, \infty[$  and for all  $k$ .

Notice that the  $k$ -dimension margins of (10) are all the same, and this for any value of  $1 \leq k \leq n$ . If  $\underline{X}$  is a *copula stationary frv* then expression (9) can be written :

$$P[\mathcal{A}_n(u)] = \psi \left( \sum_{i=1}^n \phi(G[t_i^n; u]) \right) \tag{11}$$

Table 1[Nelsen, 1999] shows three important Archimedean generators for copulas. The distribution (11) with the Clayton generator was already used for clustering of functional data coming from the symbolic data analysis framework (see [Vrac *et al.*, 2001] and [Cuvelier and Noirhomme-Fraiture, 2005]). Unfortunately the above limit is almost always null for Archimedean copulas when  $n \rightarrow \infty$  (see [Cuvelier and Noirhomme-Fraiture, 2007])!

**Proposition 1.** *If for  $u \in \mathcal{L}_2(D) : G(t; u) < 1, \forall t \in \mathcal{D}$ , then*

$$\lim_{q \rightarrow \infty} \psi \left[ \sum_{i=1}^q \phi(G[t_i^n; u]) \right] = 0 \tag{12}$$

Another objection to the use of this type of joint distribution is something which we could call *volumetric behavior*.

**Definition 11.** A function  $u \in \mathcal{L}_2(D)$  is called a *functional quantile* of value  $p$ , written  $Q_p$ , if

$$G(t; Q_p) = p, \forall t \in \mathcal{D} \tag{13}$$

The functional quantile  $Q_p$  can be seen as the level curve of value  $p$ . Now let us remark that for a functional quantile :

$$P[\mathcal{A}_n(Q_p)] = \psi \left[ \sum_{i=1}^n \phi(G[t_i^n; Q_p]) \right] = \psi(n \cdot \phi(p)) < p \tag{14}$$

And it is easy to see that, if  $n < m$  then  $\psi[m\phi(p)] < \psi[n\phi(p)]$ , and thus, the more we try to have a better approximation for a *functional quantile* of value  $p$ , the more we move away from reference value  $p$  toward zero. A simple way to avoid these two problems is to use the notion of quasi-arithmetic mean, concept which was studied by [Kolmogorov, 1930], [Nagumo, 1930] and [Aczel, 1966].

**Definition 12.** Let  $[a, b]$  be a closed real interval, and  $n \in \mathbb{N}_0$ . A quasi-arithmetic mean is a function  $M : [a, b]^n \rightarrow [a, b]$  defined as follows:

$$M(x_1, \dots, x_n) = \psi \left( \frac{1}{n} \sum_{i=1}^n \phi(x_i) \right) \quad (15)$$

where  $\phi$  is a continuous strictly monotonic real function.

We show below that if we use the generator for Archimedean copulas in (15), we define a cumulative distribution function built from one-dimensional distributions. It's easy to prove the following lemma.

**Lemma 1.** Let  $n \in \mathbb{N}_0$ ,  $F$  be a one dimensional cdf, and  $\phi$  a generator of Archimedean copula, then

$$F^*(x) = \psi \left( \frac{1}{n} \cdot \phi(F(x)) \right) \quad (16)$$

is also a cdf.

In various situations one can apply increasing transformations to the data without destroying the underlying dependence structure. This is classical in multivariate extreme value theory. And for these kind of transformation the copulas does not change.

**Proposition 2.** Let  $n \in \mathbb{N}_0$ ,  $\{F_i | 1 \leq i \leq n\}$  be a set of one dimensional cdf, and  $\phi$  a generator of Archimedean copula, then

$$H(x_1, \dots, x_n) = \psi \left( \frac{1}{n} \sum_{i=1}^n \phi(F_i(x_i)) \right) \quad (17)$$

is a multivariate cdf.

*Proof.* By the above lemma we have that the functions  $F_i^*(x)$  are cdf, and as  $\phi$  is an ‘‘Archimedean generator’’ so expression (10) is a copula, and thus  $\psi(\sum_{i=1}^n \phi(F_i^*(x_i)))$  is a multivariate cdf.  $\square$

We call the distributions given by the expression (17) the *Quasi-Arithmetic Mean of Margins (QAMM)* distributions. Now if we use a QAMM distribution in expression (5) :

$$P[\mathcal{A}_n(u)] = \psi \left[ \frac{1}{|\mathcal{D}|} \sum_{i=1}^n \Delta_i \cdot \phi(G[t_i^n; u]) \right] \quad (18)$$

then for each  $n \in \mathbb{N}$  we have an approximation, and the limit of the above expression is not always null.

**Definition 13.** Let  $\underline{X}$  be a frv,  $u \in \mathcal{L}_2(\mathcal{D})$ ,  $G$  its *Surface of Distributions* and  $\phi$  a generator of Archimedean Copulas. We define the *Quasi-Arithmetic Mean of Margins Limit (QAMML)* distribution of  $\underline{X}$  by :

$$F_{\underline{X}, \mathcal{D}}(u) = \lim_{n \rightarrow \infty} P[\mathcal{A}_n(u)] = \psi \left[ \frac{1}{|\mathcal{D}|} \cdot \int_{\mathcal{D}} \phi(G[t; u]) dt \right] \quad (19)$$

In fact transformation (16) can be seen like giving an importance to the margins in proportion with the length of an interval  $[t_i^n, t_{i+1}^n]$  in the approximation of  $F_{\underline{X}, \mathcal{D}}(u)$ .

## 4 QAMML properties

First it is easy to see that the *QAMML* distribution preserves the *functional quantiles*.

**Proposition 3.** *If  $Q_p \in \mathcal{L}_2(\mathcal{D})$  is a functional quantile of value  $p$ , then  $F_{\underline{X}, \mathcal{D}}(Q_p) = p$*

Now, what is the difference between the quasi-arithmetic mean of margins and the classical mean? Let

$$p = \frac{1}{\|\mathcal{D}\|} \int_{\mathcal{D}} G[t; u] dt \quad (20)$$

and let us define the function  $\epsilon_p[t; u] = G[t; u] - p$ . Thus we can use the following Taylor's approximation for all  $t$  (recall that  $0 \leq p, G[t; u] \leq 1$ ):

$$\begin{aligned} \phi(G[t; u]) &= \phi(p + \epsilon_p(t; u)) \\ &= \phi(p) + \phi'(p)\epsilon_p[t; u] + \phi''(p)\frac{\epsilon_p^2[t; u]}{2} + o(\epsilon_p^2[t; u]) \end{aligned} \quad (21)$$

and then

$$\begin{aligned} F_{\underline{X}, \mathcal{D}}(u) &= \psi \left[ \phi(p) + \frac{\phi'(p)}{|\mathcal{D}|} \int_{\mathcal{D}} \epsilon_p[t; u] dt + \frac{\phi''(p)}{2|\mathcal{D}|} \int_{\mathcal{D}} \epsilon_p[t; u] dt + o(\epsilon_p^2) \right] \\ &= \psi \left[ \phi(p) + \phi'(p) \mathbb{E}(\epsilon_p) + \frac{\phi''(p)}{2} \text{var}(\epsilon_p) + o(\epsilon_p^2) \right] \\ &= \psi \left[ \phi(p) + \frac{\phi''(p)}{2} \text{var}(\epsilon_p) + o(\epsilon_p^2) \right] \end{aligned} \quad (22)$$

and so like  $\psi$  is a decreasing function,  $\phi''(p) \geq 0$  and as  $\text{var}(\epsilon) \geq 0$  we can see that  $F_{\underline{X}, \mathcal{D}}(u)$  decreases with the variance of the differences between the function  $u$  and the quantile function associated to the value of the arithmetic mean of  $G$  along  $u$ . And so the *QAMML* distribution is equal to the arithmetic mean only in the case of quantile functions.

## Conclusion

In this paper we do not propose a new method in Functional Data Analysis but a new Probabilistic tool. Like in the real case, we can hope that this tool can be used for analysis of functional data, like in mixture decomposition, statistical tests,... Moreover, several ways to improve the tool exist. By example let us note that the *QAMML* definition (see (19)) uses a uniform distribution over  $\mathcal{D}$ : other distributions could be considered (see [De Finetti, 1931]).

## References

- [Aczel, 1966]J Aczel. *Lectures on Functional Equations and Their Applications*. Mathematics in Science and Engineering. Academic Press, New York and London, 1966.
- [Bartlett, 1978]M S Bartlett. *An introduction to stochastic processes*. Cambridge University Press, Cambridge, 1978.
- [Burril, 1972]C W Burril. *Measure, integration and probability*. McGraw-Hill, New-York, 1972.
- [Cox and Miller, 1965]D R Cox and HD Miller. *The theory of stochastic processes*. Methuen, London, 1965.
- [Cuvelier and Noirhomme-Fraiture, 2005]E. Cuvelier and M. Noirhomme-Fraiture. Clayton copula and mixture decomposition. In *ASMDA 2005*, pages 699–708, 2005.
- [Cuvelier and Noirhomme-Fraiture, 2007]E. Cuvelier and M. Noirhomme-Fraiture. Classification de fonctions continues l'aide d'une distribution et d'une densité définies dans un espace de dimension infinie. In *Extraction et gestion des connaissances EGC'2007*, pages 679–690, 2007.
- [De Finetti, 1931]B De Finetti. Sul concetto di media. *Giornale dell' Istituto Italiano degli Attuari*, 2:369–396, 1931.
- [Diday, 2002]E. Diday. Mixture decomposition of distributions by copulas. In *Classification, Clustering and Data Analysis*, pages 297–310, 2002.
- [Gihman and Skorohod, 1974]I I Gihman and A V Skorohod. *The theory of stochastic process*. Die grundlehren der mathematischen wissenschaften in einzel-darstellungen. Springer, Berlin, 1974.
- [Kolmogorov, 1930]A Kolmogorov. Sur la notion de moyenne. *Rendiconti Accademia dei Lincei*, 12(6):388–391, 1930.
- [Meyn and L, 1993]S P Meyn and Tweedie R L. *Markov chains and stochastic stability*. Communications and Control. Springer-Verlag, New York, 1993.
- [Nagumo, 1930]M Nagumo. Über eine klasse der mittelwerte. *Japan Journal of Mathematics*, 7:71–79, 1930.
- [Nelsen, 1999]R.B. Nelsen. *An introduction to copulas*. Springer, London, 1999.
- [Ramsay and Silverman, 2005]J O Ramsay and B W Silverman. *Functional Data Analysis*. Springer Series in Statistics. Springer, New-York, 2005.
- [Stirzaker, 2005]D Stirzaker. *Stochastic processes and models*. Oxford University Press, Oxford, 2005.
- [Vrac et al., 2001]Mathieu Vrac, Edwin Diday, Alain Chédin, and Philippe Naveau. Mélange de distributions de distributions, décomposition de mélange de copules et application à la climatologie. In *Actes du VIIIème congrès de la Société Francophone de Classification*, pages 348–355, 2001.