



HAL
open science

Space-filling designs with a Dirichlet distribution for mixture experiments

Astrid Jourdan

► **To cite this version:**

Astrid Jourdan. Space-filling designs with a Dirichlet distribution for mixture experiments. 2021. hal-03418561v1

HAL Id: hal-03418561

<https://hal.science/hal-03418561v1>

Preprint submitted on 7 Nov 2021 (v1), last revised 15 Mar 2023 (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.

Space-filling designs with a Dirichlet distribution for mixture experiments

Astrid JOURDAN

ETIS UMR 8051, CY Paris University, 95000 Cergy, France

ABSTRACT : Uniform designs are widely used for experiments with mixtures. The uniformity of the design points is usually evaluated with a discrepancy criterion. In this paper, we propose a new criterion to measure the deviation between the design point distribution and a Dirichlet distribution. The support of the Dirichlet distribution, defined by the set of d -dimensional vectors whose entries are real numbers in the interval $[0,1]$ such that the sum of the coordinates is equal to 1, is suitable for mixture experiments. Depending on its parameters, the Dirichlet distribution allows symmetric or asymmetric, uniform or more concentrated point distribution. The difference between the empirical and the target distributions is evaluated with the Kullback-Leibler divergence. We use two methods to estimate the divergence: the plug-in estimate and the nearest-neighbor estimate. The resulting two criteria are used to build space-filling designs for mixture experiments. They are also compared with two existing criteria in the case of uniform designs.

1. Introduction

Mixture experiments consist in varying the proportions of some components involved in a physico-chemical phenomenon, and observe the resulting change on the response. The proportions of the mixture components vary between 0 and 1 and they must sum to 1.0 for each run in the experiment. The experimental region is reduced to a $(d-1)$ -dimensional simplex,

$$S^{d-1} = \{(x_1, \dots, x_d) \mid x_1 + \dots + x_d = 1, x_k \geq 0\},$$

where x_k is the proportion of the k th component, $k=1, \dots, d$.

The purpose of design for mixture experiments is to define a set of points in the simplex to catch as much information about the response as possible. Since Scheffé (1958) many authors have investigated designs for mixture experiments. The pioneers (Scheffé, 1958, Kiefer, 1961, Cornell, 1981), defined optimal designs for linear and quadratic mixture models. An alternative approach of model-free designs is proposed by Wang and Fang (1990) and Fang and Wang (1994). The goal is to uniformly cover the experimental region. The main idea is to generate a uniform design on the $(d-1)$ dimensional unit cube as explained in Hickernell (1998) or in Fang *et al.* (2005). Then they apply a mapping function to put the points in the simplex S^{d-1} . Following this principle, many articles suggested improvements specially to take into account complex constraints on the components, Fang and Yang (2000), Prescott (2008), Borkowski and Piepe (2009), Ning *et al.* (2011), Liu and Liu (2016).

The former design in the unit cube is uniform in the sense that the points minimize a discrepancy criterion. The discrepancy measures the distance between uniform distribution and the empirical distribution of the design points. It is not guarantee to conserve the uniformity after the mapping function. Some authors defined criteria to assess the uniformity of design for mixture experiments. Fang and Wang (1994) proposed to use the mean square distance (MSD), Borkowski and Piepe (2009) suggested the root mean squared distance, the maximum distance and the average distance, Chuang and Hung (2010) defined the central composite discrepancy. All these criteria require to compute the distance between the design points and the points of a much larger uniform set of points. The computational cost limits their usefulness in practice. To avoid this problem, Ning *et al.* (2011) generalized the star discrepancy and proposed a new discrepancy, DM2_Discrepancy, to measure the uniformity of designs for mixtures. They also gave a computational formula of the DM2_discrepancy

only based on the design points, which is useful in practice, specially to use it in an optimization algorithm to build a uniform design for mixture experiments.

In the same way, we defined in this paper a new criterion to measure the distribution of the design points in the simplex S^{d-1} . The purpose is to obtain uniform designs, and more generally designs with a Dirichlet distribution. Depending on its parameters, the Dirichlet distribution allows to obtain symmetric and asymmetric distributions, designs with points uniformly spread in the simplex or more concentrated in the center¹. We used the Kullback-leibler (KL) divergence to measure the difference between the design point distribution and the Dirichlet distribution. The criterion is an estimator of the KL divergence computed with the design points. The KL divergence has already been used to define space-filling criteria but for an hypercube experimental domain (Jourdan and Franco, 2009, 2010). The target distribution was the uniform distribution on the unit hypercube and the criterion was reduced to the estimation of the Shannon entropy. In this paper, we adapt the criterion to the Dirichlet distribution.

In section 2, we define the criterion from the Kullback-leibler divergence and the Dirichlet distribution. In section 3, we propose two methods to estimate the criterion. In section 4, we give some numerical examples.

2. Design points with a Dirichlet distribution

Suppose that the design points x_1, \dots, x_n , are n independent observations of the random vector $X=(X_1, \dots, X_d)$ with absolutely continuous density function f concentrated on the simplex S^{d-1} . The aim is to select the design points in such a way as to have the corresponding empirical distribution “close” to the Dirichlet distribution.

Dirichlet distribution is a family of continuous multivariate probability distributions parameterized by a vector α of positive reals. The support of the Dirichlet distribution is the $(d-1)$ -simplex S^{d-1} . Its probability density function is

$$g(x) = \frac{1}{B(\alpha)} \prod_{k=1}^d (x_k)^{\alpha_k - 1},$$

where x belongs to the $(d-1)$ -simplex S^{d-1} , $\alpha=(\alpha_1, \dots, \alpha_d)$ with $\alpha_i > 0$, and $B(\alpha)$ is the normalizing constant,

$$B(\alpha) = \frac{\prod_{k=1}^d \Gamma(\alpha_k)}{\Gamma(\alpha_0)} \quad \text{with } \alpha_0 = \sum_{k=1}^d \alpha_k.$$

A common special case is the symmetric Dirichlet distribution, where all of the elements making up the parameter vector α have the same value α , called the concentration parameter. When $\alpha=1$, the symmetric Dirichlet distribution is equivalent to a uniform distribution over the $(d-1)$ -simplex S^{d-1} . It is called the flat Dirichlet distribution.

The aim is to generate n points in the simplex with a distribution as close as possible of a Dirichlet distribution. On Figures 1 and 2 (starting design), we can see that a simple random generation of the Dirichlet distribution is not efficient to obtain a good point distribution. Especially in the case of the flat Dirichlet distribution (Figures 1a and 2d), the points do not uniformly cover the simplex: some points are very close to each other while some areas are empty.

We defined a criterion to measure the “distance” between the point distribution and the Dirichlet distribution. The criterion is then used in an optimization algorithm to build a set of points with the expected distribution.

¹ Plots of probability density functions from the Dirichlet distribution : <https://en.wikipedia.org/wiki/File:Dirichlet.pdf>

There are different ways to measure the difference between two distributions. In this paper, we use the Kullback-Leibler divergence to evaluate the deviation between two probability density functions f and g ,

$$I(f,g) = \int_{S^{d-1}} f(x) \log\left(\frac{f(x)}{g(x)}\right) dx .$$

In this paper g is the density function of the Dirichlet distribution and f is the density function of the design points. This integral can be written as the expected value of a random vector X with Dirichlet distribution,

$$I(f,g) = E\left[\log\left(\frac{f(X)}{g(X)}\right)\right].$$

If we consider that the design points $D=\{x_1,\dots,x_n\}$ are n i.i.d. realizations of a Dirichlet distribution, the Monte Carlo method gives an unbiased and consistent estimator,

$$\hat{I}(f,g) = \frac{1}{n} \sum_{i=1}^n \log\left(\frac{f(x_i)}{g(x_i)}\right) \quad (1)$$

This expression is not a computational formula since the density function f is unknown. There are two common ways to estimate integral I : the plug-in estimate which consists in replacing the density function f by its kernel estimate, and the nearest-neighbor estimate.

The two estimations are no more unbiased. However, having a bias is not a problem in our application, if the bias is fixed for a given n and d . The goal is not to obtain an accurate estimate of the integral but a criterion to compare two set of points in the optimization algorithm. We say that a design D_1 is better than a design D_2 if

$$\hat{I}(f_1,g) \leq \hat{I}(f_2,g)$$

with f_1 and f_2 the density functions associated to D_1 and D_2 respectively.

The optimization algorithm is an adaptation of the exchange algorithm described in Jin *et al.* (2005)

Exchange Algorithm

1. Simulate n points from a Dirichlet distribution²
2. Randomly select a new point in the simplex S^{d-1}
3. For i in 1 to n
 - Build a new design by replacing the i^{th} point by the new point
 - Compute the criterion value of the new design
 - Replace the current design with the new one iff the exchange improves the criterion
- End for i
4. Repeat steps 2 and 3 until terminating condition is met

This algorithm is a simple local search algorithm and could quickly find a locally optimal design (more sophisticated algorithms are given in Fang *et al.*, 2005, Chapter 4, §4.2). The resulting design depends more or less on the initial setting. Hence, once should repeat the algorithm with several different random starting designs and select the best design among the generated designs. However, in Section 4, several designs are built from only one initialization in order to study the influence of the starting design.

² Let y_1,\dots,y_d be i.i.d realizations of the Gamma distribution with $y_i \sim \Gamma(\alpha_i,1)$. The random vector $x=y/S$ where $S=y_1+\dots+y_d$ has a Dirichlet distribution with parameter $\alpha=(\alpha_1,\dots,\alpha_d)$.

3. Estimation of the criterion

In this section we propose two methods to estimate the unknown density function f in (1). In each case we explain our choices (kernel, bandwidth, k in the k -nearest neighbor distance) and we give a computational formula for the criterion.

3.1. Plug-in estimate

The unknown density function f is estimated with the design points $D=\{x_1, \dots, x_n\}$ by a kernel method (Scott, 1992)

$$\hat{f}(x) = \frac{1}{n|H|^{1/2}} \sum_{i=1}^n K(H^{-1/2}(x - x_i))$$

where K is a multivariate kernel and H is the bandwidth matrix (symmetric and positive definite matrix).

Choice of the bandwidth

The choice of the bandwidth matrix has a great influence on the accuracy of the estimation. Joe (1989) shows that in the case where f is estimated by a kernel method, the bias in the estimation of $I(f, g)$ depends on the sample size n , the dimension d , and the bandwidth matrix H . When constructing an optimal design, the size n and the dimension d are fixed. The bandwidth still needs to be fixed so that the bias does not vary during the optimization algorithm.

Usually the bandwidth matrix is chosen to be proportional to the covariance matrix of the data. This solution implies that H varies during the optimization algorithm. An idea to fix it, is to replace the covariance matrix of the data by the target covariance matrix, *i.e* the covariance matrix of the Dirichlet distribution. Unfortunately, this matrix is singular. Then, even if the variables are correlated, we simplify the bandwidth matrix into a diagonal matrix with the Scott's rule (1992), $H = \text{diag}(h_1^2, \dots, h_d^2)$ with $h_j = n^{-1/(d+4)} \hat{\sigma}_j$. As previously, the estimation of the standard deviation of the j^{th} component, $\hat{\sigma}_j$ changes at each iteration of the optimization algorithm. In order to fix the bias, we replace it with the standard deviation of the target distribution,

$$\hat{\sigma}_j = \frac{1}{\alpha_0} \sqrt{\frac{\alpha_j(\alpha_0 - \alpha_j)}{(\alpha_0 + 1)}}.$$

Choice of the kernel

It is known that the shape of the kernel has a minor influence on the estimation (Silverman, 1986). We have chosen to use a multidimensional Gaussian kernel,

$$K(z) = (2\pi)^{-d/2} e^{-\frac{1}{2}\|z\|^2}.$$

A kernel of finite support (Epanechnikov, uniform,...) might have seemed more appropriate for the Dirichlet distribution. But, in our application, the kernel function has input values z in the interval $[0, d/h^2]$. This interval becomes very large when the size and the dimension increase, and the probability for z to be inside the kernel support then becomes very low. The estimation of f is then almost constant during the optimization process, and the criterion therefore does not allow to compare the designs.

Finally,

$$\hat{I}(\hat{f}, g) = \frac{1}{n} \sum_{i=1}^n \left[\log \left\{ \frac{1}{n} \sum_{j=1}^n \frac{(2\pi)^{-d/2}}{h_1 \dots h_d} \exp \left(-\frac{1}{2} \sum_{k=1}^d \left(\frac{x_{jk} - x_{ik}}{h_k} \right)^2 \right) \right\} - \log \left\{ \prod_{k=1}^d x_{ik}^{\alpha_k - 1} \right\} + \log(B(\alpha)) \right] \quad (2)$$

where x_{ik} is the k th component of the i th design point, $i=1, \dots, n$ and $k=1, \dots, d$, and

$$h_j = n^{-1/(d+4)} \frac{1}{\alpha_0} \sqrt{\frac{\alpha_j(\alpha_0 - \alpha_j)}{(\alpha_0 + 1)}}.$$

By removing the terms independent of the design points, we obtain a simplified criterion, especially for the symmetric Dirichlet distribution,

$$C_{\text{kern}}(D) = \sum_{i=1}^n \left[\log \left(\sum_{j=1}^n e^{-\frac{1}{2} \left\| \frac{x_{jk} - x_{ik}}{h} \right\|^2} \right) \right] - (\alpha - 1) \sum_{i=1}^n \sum_{k=1}^d \log(x_{ik}) \quad (3)$$

where

$$h = n^{-1/(d+4)} \frac{1}{d} \sqrt{\frac{d-1}{d\alpha+1}}.$$

The calculus is detailed in the appendix.

3.2. Nearest-neighbor estimate

Wang *et al.* (2006) and Leonenko *et al.* (2008) proposed to estimate the Kullback-Leibler divergence with the k-nearest neighbor density estimation.

Let $\rho(x, y)$ denote the Euclidian distance between two points x and y of \mathbb{R}^d . We note $\rho^{(1)}(x, S) \leq \rho^{(2)}(x, S) \leq \dots \leq \rho^{(m)}(x, S)$, the order distances between $x \in \mathbb{R}^d$ and $S = \{y_1, \dots, y_m\}$ a set of points of \mathbb{R}^d such that $x \notin S$. $\rho^{(k)}(x, S)$ is the k-nearest-neighbor distance from x to points of S . The previous authors demonstrated that the following estimate of $l(f, g)$ with the design points $D = \{x_1, \dots, x_n\}$ is asymptotically unbiased and consistent,

$$\hat{l}(\hat{f}, \hat{g}) = \frac{1}{n} \sum_{i=1}^n \left[-\log \left\{ (n-1) e^{-\Psi^{(k)}} V_d \left(\rho^{(k)}(x_i, D_{-i}) \right)^d \right\} - \log \left\{ \prod_{k=1}^d x_{ik}^{\alpha_k - 1} \right\} + \log(B(\alpha)) \right] \quad (4)$$

with Ψ the digamma function, V_d the volume of the unit ball in \mathbb{R}^d and $D_{-i} = D \setminus \{x_i\}$. By removing the terms independent of the design points, we obtain the following criterion for a symmetric Dirichlet distribution,

$$C_{\text{nn}}(D) = - \sum_{i=1}^n \log \left\{ \left(\rho^{(k)}(x_i, D_{-i}) \right)^d \right\} - (\alpha - 1) \sum_{i=1}^n \sum_{k=1}^d \log(x_{ik}) \quad (5)$$

The choice of k is discussed in the next section.

Remark : Note that the criteria C_{kern} and C_{nn} are reduced to their first term for the flat Dirichlet (uniform) distribution ($\alpha=1$).

Remark: The complexity of the two criteria is $O(n^2)$.

4. Numerical tests

In this section we built designs with the two criteria C_{kern} and C_{nn} and the optimization algorithm presented in section 2 for different values of d , n and α . For each configuration, we built 30 designs in order to take into account the random initialization. The same random designs are used at the beginning of the optimization algorithm for both criteria C_{kern} and C_{nn} .

For each resulting design we compute two existing criteria, the MSD criterion defined by Fang and Wang (1994), and the DM2 criterion proposed by Ning *et al.* (2011).

4.1. Impact of the k value

In order to study the impact of the k value for the resulting C_{nn} design, different values have been tested between 1 and $2n/3$.

Table 1 gives the squared correlation between the C_{nn} , DM2 and MSD criteria and the k value used to build the 30 designs with the C_{nn} criterion. The C_{nn} criterion and the k value are strongly correlated, the criterion decreases when k increases. It is therefore important to use the same value of k to compare the designs in the optimization process. The square correlation is very small for the two other criteria. It seems that the choice of k has no impact on the resulting design.

	d=3, n=30			d=4, n=40			d=5, n=50			d=6, n=60			d=7, n=70		
	DM2	MSD	C_{nn}	DM2	MSD	C_{nn}	DM2	MSD	C_{nn}	DM2	MSD	C_{nn}	DM2	MSD	C_{nn}
MSD	-0,48	1		-0,26	1		-0,16	1		-0,40	1		-0,25	1	
C_{nn}	0,34	-0,34	1	0,33	-0,53	1	0,43	-0,47	1	-0,12	0,32	1	0,14	-0,34	1
k	-0,04	0,20	-0,92	-0,34	0,48	-0,95	-0,36	0,41	-0,95	0,14	-0,31	-0,95	-0,02	0,33	-0,95

Table 1. Squared correlation between the C_{nn} , DM2 and MSD criteria and the k value used to build the 30 designs with the C_{nn} criterion and with $\alpha=1$.

4.2. Design comparison

Figures 1 and 2 give some examples of designs obtained with C_{kern} and C_{nn} criteria from the same starting design. The criteria work well since the design points are more evenly spread in the simplex than the original ones for $\alpha=1$. We note that C_{kern} tends to put the points on the boundaries while the points obtained with C_{nn} are more concentrated in the middle. If we constraint the algorithm to put points near the boundaries, it degrades the C_{nn} criterion. Inversely, more points in the middle degrades the C_{kern} criterion.

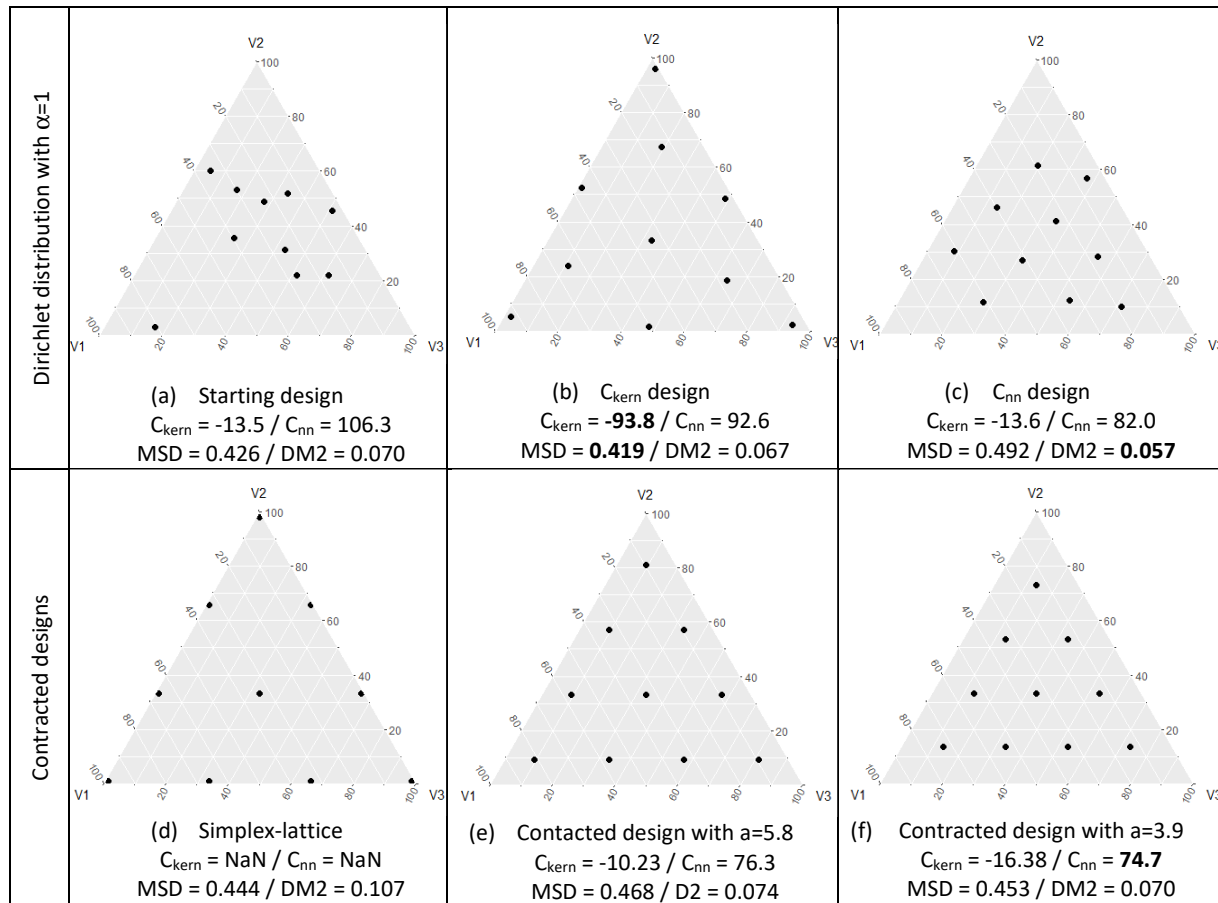


Figure 1. Uniform designs for mixture experiments with $d=3$ and $n=10$. The starting design points are n i.i.d random generation of the Dirichlet distribution (a). (b) and (c) are the resulting designs of the optimization algorithm with the same initialization (a). (e) and (f) are the contracted form of the simplex-lattice (d) (§4.2).

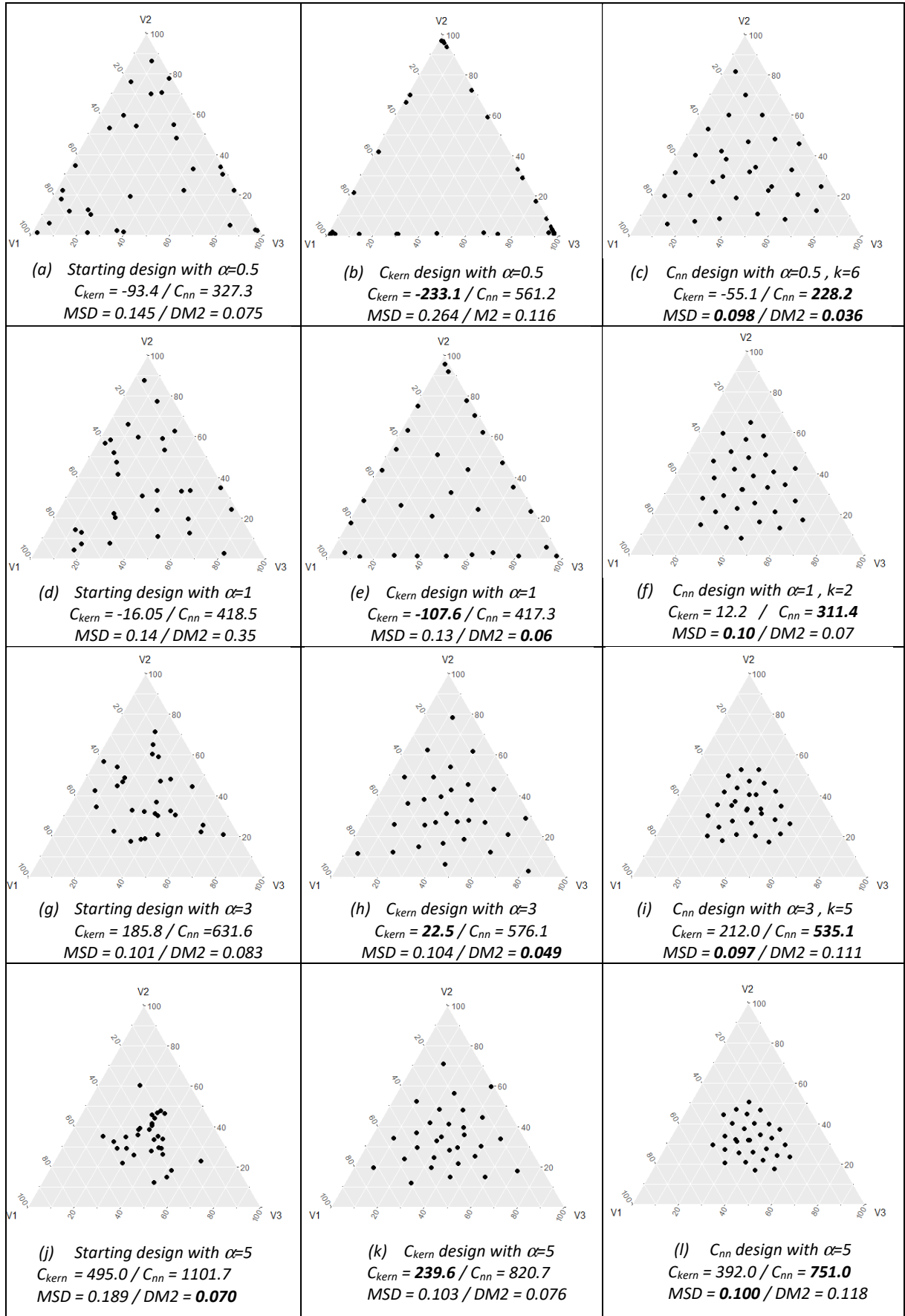


Figure 2. Evolution of the designs with for $d=3$ and $n=30$ built with C_{kern} or C_{nn} criterion when α increases. The starting design points are n.i.d random generation for the Dirichlet distribution (left column). The middle and right columns are the resulting designs of the optimization algorithm with the same initialization given in the left column.

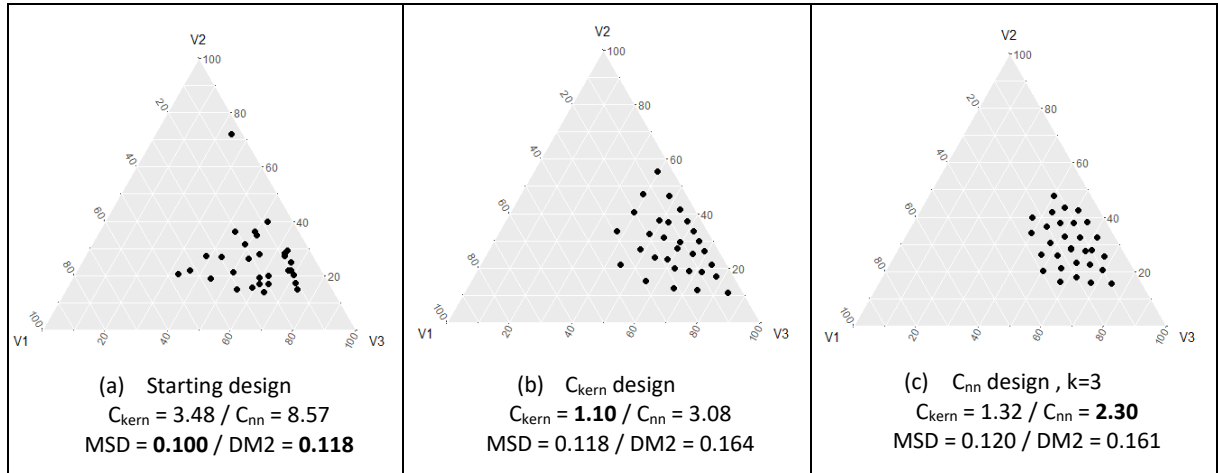


Figure 3. Designs built with C_{kern} or C_{nn} criterion with $d=3$ and $n=10$ for an asymmetric Dirichlet distribution with $\alpha=(2,4,8)$. The starting design points are *n.i.d* random generation of the Dirichlet distribution.

The boxplots in Figures 4 and 5 represent the 30 values of DM2 and MSD criteria for the C_{kern} and C_{nn} designs. For the C_{nn} designs, we notice on Figure 2 that the points are concentrated in the center of the simplex when α increases. The DM2 criterion reflects this behavior well since it increases with α (Figure 4). Its minimum value is reached for $\alpha=0.5$. This indicates that the C_{nn} design have a uniform distribution for $\alpha=0.5$ (in the sense of the DM2 criterion). This result does not match with the definition of the Dirichlet distribution since the uniform distribution is reached with $\alpha=1$. However, we can see graphically on Figure 2c that the points seem indeed more uniformly distributed in the simplex with $\alpha=0.5$. On the contrary, the MSD criterion does not appear to be a good indicator of uniform point distribution. It increases when α decreases (Figure 5). Its minimum value is reached for $\alpha=5$. However, we can see on Figure 2l that in this case the design points are very concentrated in the center of the simplex. The same behavior of the MSD criterion is observed for the C_{kern} designs.

For C_{kern} designs, the DM2 criterion decreases with α until $\alpha=3$ (Figure 4). Graphically we can see that for $\alpha=0.5$ or $\alpha=1$ the points are on the edges (Figures 2b and 2e). Then the DM2 value increases for $\alpha=5$, i.e. when the design points are concentrated in the middle (Figure 2k). Finally the best value is reached for $\alpha=3$ and not for $\alpha=1$.

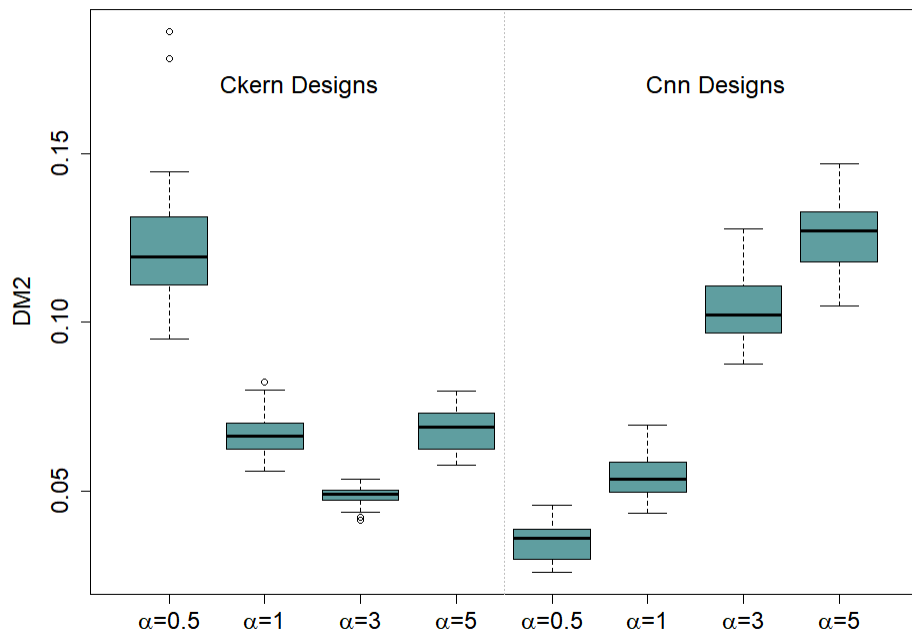


Figure 4. Boxplot of the DM2 criterion values for the sampling of 30 designs with $d=3$ and $n=30$ with $\alpha=0.5, 1, 3$ and 5 . Designs built with C_{kern} criterion are on the left and those built with C_{nn} are on the right.

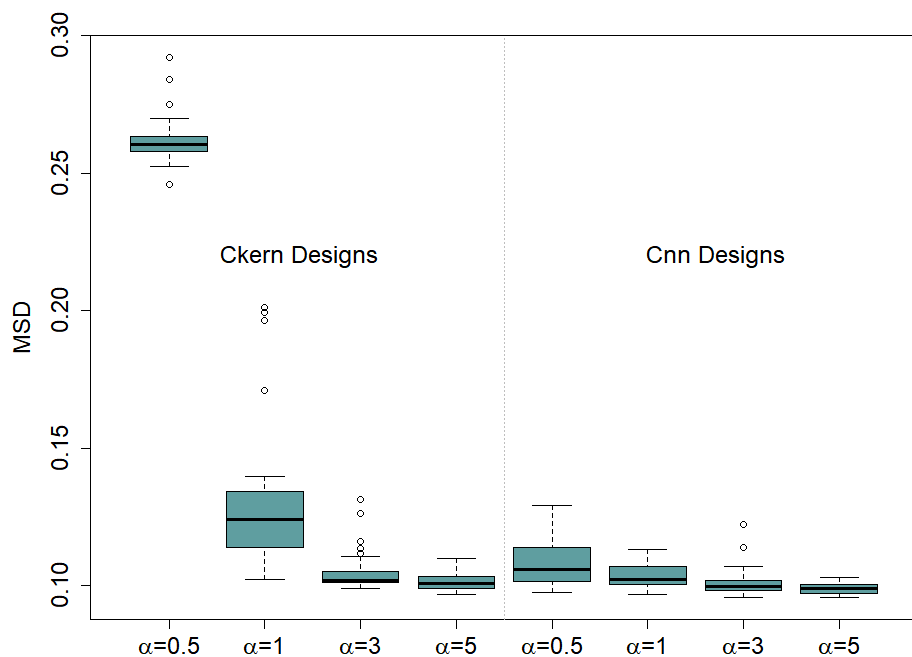


Figure 5. Boxplot of the MSD criterion values for the sampling of 30 designs with $d=3$ and $n=30$ and with $\alpha=0.5, 1, 3$ and 5 . Designs built with C_{kern} criterion are on the left and those built with C_{nn} are on the right.

4.3. Concentrated design

An alternative to build uniform design for mixture experiments is the contraction of a simplex-lattice (Scheffé, 1958). The points of a simplex-lattice seem to be uniformly distributed on S^{d-1} but most of

them lies on the boundaries (Figure 1d). The mixture is then reduced to d-1 or d-2 ingredients and one or two ingredients are not in the mixture. Fang and Wang (1994) proposed to keep the simplex-lattice pattern while moving the points towards the centroid of the simplex. An example of a lattice-simplex and the contracted design is given in Table 2. The smaller the contraction constant a , the more the points are concentrated in the center (Figures 1e and 1f). Prescott (2008) and Ning *et al.* (2011) used the MSD and DM2 criteria respectively to find the best value of a . In the same way, we optimize the C_{kern} and C_{nn} criteria to determine a (Figure 6).

The C_{kern} criterion does not converge. It decreases when a increases. This confirms that the C_{kern} criterion tends to push the points on the edges of the simplex. On the contrary, the best value obtained with the MSD criterion is small ($a=3.9$). This also confirms that the MSD criterion favors the points in the center. The C_{nn} and DM2 criteria have the same behavior and give nearly the same value ($a=5.8$ with C_{nn} and $a=5.26$ with DM2)

Simplex-lattice design			Contracted design		
X ₁	X ₂	X ₃	X ₁	X ₂	X ₃
1	0	0	$1-1/a$	$1/(2a)$	$1/(2a)$
0	1	0	$1/(2a)$	$1-1/a$	$1/(2a)$
0	0	1	$1/(2a)$	$1/(2a)$	$1-1/a$
$2/3$	$1/3$	0	$2/3-1/(2a)$	$1/3$	$1/(2a)$
$1/3$	$2/3$	0	$1/3$	$2/3-1/(2a)$	$1/(2a)$
$2/3$	0	$1/3$	$2/3-1/(2a)$	$1/(2a)$	$1/3$
$1/3$	0	$2/3$	$1/3$	$1/(2a)$	$2/3-1/(2a)$
0	$2/3$	$1/3$	$1/(2a)$	$2/3-1/(2a)$	$1/3$
0	$1/3$	$2/3$	$1/(2a)$	$1/3$	$2/3-1/(2a)$
$1/3$	$1/3$	$1/3$	$1/3$	$1/3$	$1/3$

Table 2. {3,3}-simplex lattice design and its contracted design.

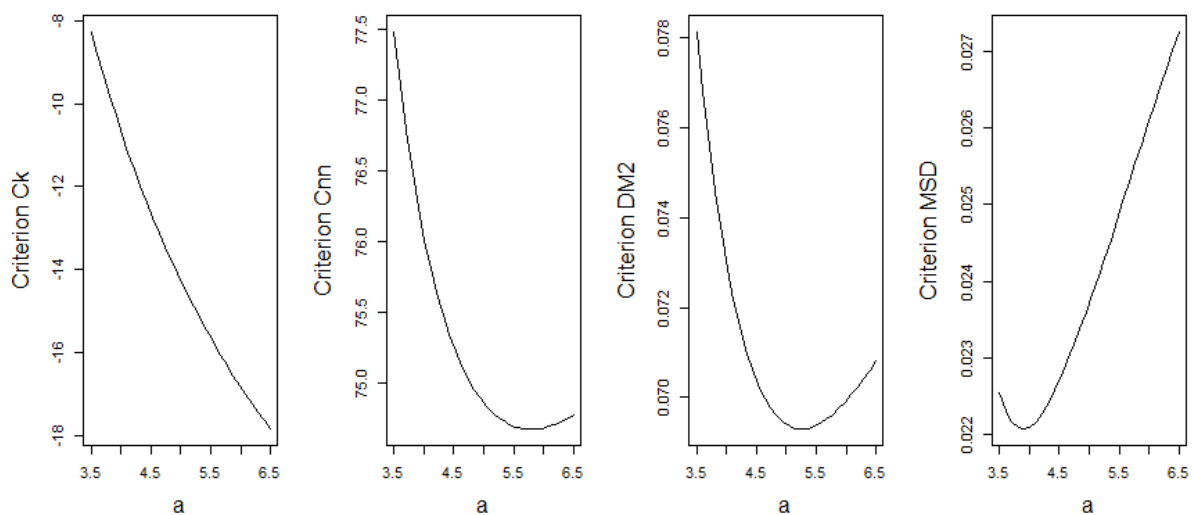


Figure 6. C_{kern} , C_{nn} , DM2 and MSD criteria against the contraction constant a . Criterion C_{nn} is computed with $k=1$. Best values of a are 5.8 with C_{nn} , 5.26 with DM2 and 3.9 with MSD.

4.4. Marginal distributions

The marginal distributions of the Dirichlet distribution are Beta distributions, $\text{Beta}(\alpha_i, \alpha_i - \alpha_0)$. In the special case of the uniform distribution ($\alpha=1$), the distributions are $\text{Beta}(1, d-1)$. The asymmetric shape of the density function implies that the small values (proportions) of the components are over-represented while the larger values are under-represented. The same behavior is observed with the contracted design (Figure 7). There is no reason to make more experiments with small values. We would also like to have a uniform distribution for each of the components. But the two objectives, uniform distribution on the simplex and uniform distribution for each component, are conflicting.

Figure 7 illustrates the component distributions with the designs of Figure 1. Figures 7a and 7b show the three distributions of X_1 , X_2 and X_3 for the C_{kern} and C_{nn} designs. Figure 7c displays only X_1 distribution of the contracted design but for different values of the contraction constant a . We can see that the marginal distributions are more uniform with the design computed with the C_{kern} criterion since they are very asymmetric for C_{nn} and contracted designs.

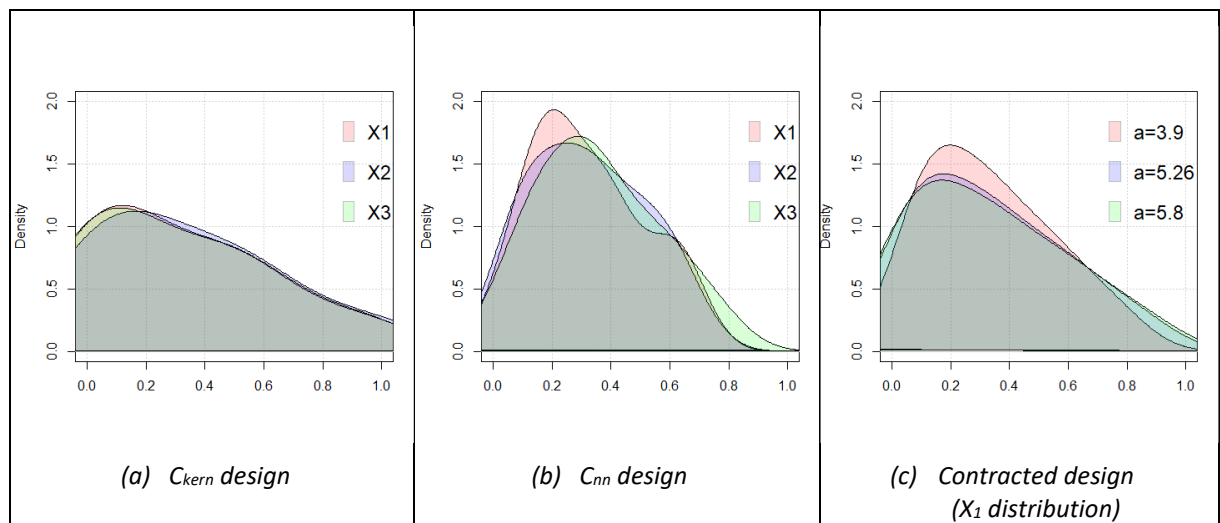


Figure 7. Component distributions for designs with $d=3$ and $n=10$ given in Figure 1.

Conclusion

In this paper we proposed a new class of designs for mixture experiments. The Dirichlet distribution allows to build design points with symmetric or asymmetric distribution, uniform or contracted distribution. The Kullback-Leibler divergence is used to measure the difference between the Dirichlet and design points distributions. We used the plugin estimate with a Gaussian kernel and the nearest neighbor estimate of the Kullback-Leibler divergence to define two criteria to assess the design point distribution. The two criteria are simplified to be used in an optimization process to build designs for mixture experiments with a target Dirichlet distribution.

The numerical tests and the comparison with two existing criteria (MSD and DM2) show that the C_{kern} and C_{nn} criteria perform well to build uniform designs. With the same α value, the C_{kern} criterion tends to push the points on the boundaries while the C_{nn} criterion concentrates the points in the middle of the simplex. The advantage of the criteria proposed in this article is that they allow the construction of designs with distributions other than uniform, symmetric or asymmetric.

However, our method presents a shortcoming. The optimization algorithm converges very slowly and requires many iterations until terminating condition is met. An early stopping of the optimization process may produce poor quality designs. One idea to speed up convergence is to choose the new

point in a large NT-net instead of choosing it randomly in the simplex (step 2 in the exchange algorithm).

Another drawback of uniform design for mixture experiments (not only C_{kern} and C_{nn} designs) is the asymmetric distribution of each component. Having a uniform distribution on the simplex S^{d-1} and a symmetric distribution on each axis seem to be two conflicting objectives. A multi-objective optimization algorithm (instead of the exchange algorithm) would allow to manage this problem. The first objective function would be one of the two criteria defined in this paper. The second objective function could be defined in order to measure the difference between the distribution of each component and a symmetric distribution with support $[0,1]$ (e.g. symmetric triangular or truncated normal distribution). As we did in this paper, the Kullback-Leibler divergence and its estimates could be used to define this second objective function.

APPENDIX

In this appendix, we give the detailed calculus to obtain formulas (2) and (3).

If we consider that the design points $D=\{x_1, \dots, x_n\}$ are n i.i.d. realizations of a Dirichlet distribution, the Monte Carlo method gives an unbiased and consistent estimator,

$$\hat{\imath}(f, g) = \frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(x_i)}{g(x_i)} \right) = \frac{1}{n} \sum_{i=1}^n \log(f(x_i)) - \log(g(x_i)).$$

By replacing g with the Dirichlet density function, we obtain,

$$\hat{\imath}(f, g) = \frac{1}{n} \sum_{i=1}^n \left[\log(f(x_i)) - \log \left\{ \prod_{k=1}^d x_{ik}^{\alpha_k - 1} \right\} + \log(B(\boldsymbol{\alpha})) \right] \quad (\text{i})$$

Note that x_{ik} must be non-zero.

If the Dirichlet distribution is symmetric, then

$$\hat{\imath}(f, g) = \frac{1}{n} \sum_{i=1}^n \log(f(x_i)) - \frac{(\alpha - 1)}{n} \sum_{i=1}^n \sum_{k=1}^d \log(x_{ik}) + \log(B(\boldsymbol{\alpha})) \quad (\text{ii})$$

The unknown density function f is estimated with the design points $D=\{x_1, \dots, x_n\}$ by a kernel method

$$\hat{f}(x) = \frac{1}{n |H|^{1/2}} \sum_{j=1}^n K(H^{-1/2}(x - x_j))$$

where $H = \text{diag}(h_1^2, \dots, h_d^2)$ and $K(z) = (2\pi)^{-d/2} e^{-\frac{1}{2}\|z\|^2}$. Then

$$\hat{f}(x) = \frac{1}{n h_1 \dots h_d} \sum_{j=1}^n K \left(\frac{x_1 - x_{j1}}{h_1}, \dots, \frac{x_d - x_{jd}}{h_d} \right) = \frac{1}{n h_1 \dots h_d} \sum_{j=1}^n (2\pi)^{-d/2} e^{-\frac{1}{2} \sum_{k=1}^d \left(\frac{x_k - x_{jk}}{h_k} \right)^2}.$$

Formula (2) is obtained by replacing $f(x_i)$ by $\hat{f}(x_i)$ in Formula (i),

$$\hat{\imath}(\hat{f}, g) = \frac{1}{n} \sum_{i=1}^n \left[\log \left(\frac{1}{n h_1 \dots h_d} \sum_{j=1}^n (2\pi)^{-d/2} e^{-\frac{1}{2} \sum_{k=1}^d \left(\frac{x_{ik} - x_{jk}}{h_k} \right)^2} \right) - \log \left\{ \prod_{k=1}^d x_{ik}^{\alpha_k - 1} \right\} + \log(B(\boldsymbol{\alpha})) \right].$$

In the symmetric case, $\alpha_1 = \dots = \alpha_d = \alpha$, then $\alpha_0 = d\alpha$ and the bandwidth is constant,

$$h_1 = \dots = h_d = h = n^{-1/(d+4)} \frac{1}{d} \sqrt{\frac{d-1}{d\alpha+1}}.$$

Formula (ii) gives

$$\begin{aligned} \hat{l}(\hat{f}, g) &= \frac{1}{n} \sum_{i=1}^n \log \left(\frac{1}{n h^d} \sum_{j=1}^n (2\pi)^{-d/2} e^{-\frac{1}{2} \sum_{k=1}^d \left(\frac{x_{ik} - x_{jk}}{h} \right)^2} \right) - \frac{(\alpha-1)}{n} \sum_{i=1}^n \sum_{k=1}^d \log(x_{ik}) + \log(B(\alpha)) \\ &= \log \left(\frac{(2\pi)^{-d/2}}{n h^d} \right) + \frac{1}{n} \sum_{i=1}^n \log \left(\sum_{j=1}^n e^{-\frac{1}{2} \sum_{k=1}^d \left(\frac{x_{ik} - x_{jk}}{h} \right)^2} \right) - \frac{(\alpha-1)}{n} \sum_{i=1}^n \sum_{k=1}^d \log(x_{ik}) + \log(B(\alpha)) \end{aligned}$$

If we remove the terms independent of the design points, we have Formula (3)

$$C_{\text{kernel}}(D) = \sum_{i=1}^n \log \left(\sum_{j=1}^n e^{-\frac{1}{2} \sum_{k=1}^d \left(\frac{x_{ik} - x_{jk}}{h} \right)^2} \right) - (\alpha-1) \sum_{i=1}^n \sum_{k=1}^d \log(x_{ik})$$

References

- Borkowski, J.J., Piepel, G.F. (2009). Uniform designs for highly constrained mixture experiments. *Journal of Quality Technology*, 41 (1), 35–47.
- Chuang S.C., Hung Y.C. (2010). Uniform design over general input domains with applications to target region estimation in computer experiments. *Computational Statistics & Data Analysis*, 54 (1), 219-232.
- Cornell, J. A. (1981). *Experiments with Mixtures, designs, models, and the analysis of mixture data*. Wiley, New York.
- Fang, K.T., Li, R., Sudjianto, A. (2005). *Design Modeling for Computer Experiments*. Chapman & Hall/CRC Press, London.
- Fang, K.T., Wang, Y. (1994). *Number-theoretic Methods in Statistics*. Chapman & Hall, London.
- Fang, K.T., Yang, Z.H. (2000). On uniform design of experiments with restricted mixture and generation of uniform distribution on some domains. *Statistics & Probability Letters*, 46, 113–120.
- Hickernell, F.J. (1998). A generalized discrepancy and quadrature error bound. *Mathematics of Computation*, 67, 299-322.
- Jin, R., Chen, W., Sudjianto, A. (2005). An efficient algorithm for constructing optimal design of computer experiments. *Journal of Statistical Planning and Inference*, 134, 268–287.
- Jourdan A. et Franco J. (2009). Plans d'expériences numériques d'information de Kullback-Leibler minimale. *Journal de la Société Française de Statistique*, 150 (2), 52-64.
- Jourdan A. et Franco J. (2010). Optimal Latin hypercube designs for the Kullback-Leibler criterion. *ASTA Advances in Statistical Analysis*, 94 (4), 341-351.
- Kiefer, J. (1961). Optimal designs for regression model, II. *The Annals of Mathematical Statistics*, 32, 298–325.
- Leonenko N, Pronzato L, Savani V. (2008). A class of Rényi information estimators for multidimensional densities. *The annals of Statistics*, 36 (5), 2153-2182.

- Liu, Y., & Liu, M. (2016). Construction of uniform designs for mixture experiments with complex constraints. *Communications in Statistics - Theory and Methods*, 45, 2172 - 2180.
- Ning J.H., Zhou Y.D., Fang K.T. (2011). Discrepancy for uniform design of experiments with mixtures, *Journal of Statistical Planning and Inference*, 141 (4), 1487-1496.
- Prescott, P. (2008). Nearly uniform designs for mixture experiments. *Communication in Statistics - Theory and Methods*, 37, 2095–2115.
- Scheffé, H. (1958). Experiments with Mixtures. *Journal of the Royal Statistical Society: Series B*, 20 (2), 344-360.
- Scott D.W. (1992). *Multivariate Density Estimation: Theory, practice and visualization*. John Wiley & Sons, New York, Chichester.
- Wang Q., Kulkarni R., Verdu S. (2006). A Nearest-Neighbor Approach to Estimating Divergence between Continuous Random Vectors. 2006 IEEE International Symposium on Information Theory, 242-246.
- Wang Y., Fang K.T. (1990). Number theoretic methods in applied statistics (II). *Chinese Annals of Mathematics - Series B*, 11,384-394.