

Random vectors

Let (Ω, \mathcal{B}, P) be a probability space. A real-valued random vector

$$\mathbf{X}(\omega) = (X_1(\omega), \dots, X_n(\omega))$$

is a measurable map from the sample space Ω into \mathbb{R}^n , i.e.,

$$\mathbf{X} : \Omega \rightarrow \mathbb{R}^n. \quad (1)$$

Each component of $\mathbf{X}(\omega)$, say $X_i(\omega)$, is a real-valued random variable. As before, we can push forward the probability measure P from Borel sets of Ω to Borel sets of \mathbb{R}^n via the mapping (1), i.e.,

$$P_{\mathbf{X}}(A) = \mathcal{B}(\mathbb{R}^n) \mapsto [0, 1] \quad (2)$$

where $P_{\mathbf{X}}(A)$ is defined as¹

$$P_{\mathbf{X}}(A) = P(\{\omega \in \Omega : \mathbf{X}(\omega) \in A\}) \quad \text{for all } A \in \mathcal{B}(\mathbb{R}^n). \quad (4)$$

Joint CDFs and PDFs. The *cumulative distribution function* (CDF) of a $\mathbf{X}(\omega)$ is defined as

$$F(x_1, \dots, x_n) = P(\underbrace{\{\omega : X_1(\omega) \leq x_1\} \cap \dots \cap \{\omega : X_n(\omega) \leq x_n\}}_{\text{event in } \mathcal{B}(\Omega) \text{ defined as intersection of } n \text{ events}}). \quad (5)$$

As before, if P is absolutely continuous with respect to the Lebesgue measure $dx_1 \cdots dx_n$ then there exists a (Lebesgue integrable) *probability density function*² (PDF) $p(x_1, \dots, x_n)$ such that

$$F(x_1, \dots, x_n) = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_n} p(y_1, \dots, y_n) dy_1 \cdots dy_n. \quad (6)$$

Equivalently, we can express $p(x_1, \dots, x_n)$ as a (weak) derivative of $F(x_1, \dots, x_n)$ as

$$p(x_1, \dots, x_n) = \frac{\partial^n F(x_1, \dots, x_n)}{\partial x_1 \cdots \partial x_n}. \quad (7)$$

The multivariate distribution function F and associated probability density function p satisfy similar properties as the properties we have seen for one one random variable (see [13] for details). For instance F is non-decreasing, with range in $[0, 1]$, etc. Similarly, p is non-negative, and it allows us to compute the probability of the event

$$\{\omega \in \Omega : \mathbf{X}(\omega) \in A\}$$

as

$$P(A) = \int_A p(x_1, \dots, x_n) dx_1 \cdots dx_n. \quad (8)$$

Frequency interpretation of the joint PDF: Consider a continuous random vector $\mathbf{X}(\omega)$ (i.e., a random vector with continuous CDF) with only two components, say $X_1(\omega)$ and $X_2(\omega)$. By using equations (5) and (6) we have

$$P(\{\omega : x_1 \leq X_1(\omega) \leq x_1 + \Delta x_1\} \cap \{\omega : x_2 \leq X_2(\omega) \leq x_2 + \Delta x_2\}) \simeq p(x_1, x_2) \Delta x_1 \Delta x_2. \quad (9)$$

¹The set

$$\mathbf{X}^{-1}(A) = \{\omega \in \Omega : \mathbf{X}(\omega) \in A\} \quad (3)$$

is known as pre-image of A under the mapping $\mathbf{X}(\omega)$.

²Technically speaking, the joint probability density function $p(x_1, \dots, x_n)$ is the Radon-Nikodym derivative of the probability measure P relative to the Lebesgue measure $dx_1 \cdots dx_n$.

Let us partition the tensor product space \mathbb{R}^2 with an evenly-spaced grid of width Δx_1 (along x_1) and Δx_2 (along x_2). Suppose we observe S realizations of $\mathbf{X}(\omega) = (X_1(\omega), X_2(\omega))$, and suppose that $n_A < S$ instances satisfy the condition

$$\{x_1 \leq X_1(\omega) \leq x_1 + \Delta x_1\} \quad \text{and} \quad \{x_2 \leq X_2(\omega) \leq x_2 + \Delta x_2\}. \quad (10)$$

Then from (9) we obtain the PDF estimate

$$p(x_1, x_2) \simeq \frac{1}{\Delta x_1 \Delta x_2} \frac{n_A}{n}. \quad (11)$$

Of course there are methods other than relative frequencies to estimate PDFs from data/observations. Among them, we have kernel methods [2] (see Figure 1) and generative modeling techniques based on diffusion [14].

Marginal CFDs and PDFs. Let $\mathbf{X}(\omega) = (X_1(\omega), X_2(\omega))$ be a random vector with joint distribution function $F(x_1, x_2)$. The distribution of the random variable $X_1(\omega)$ can be obtained from $F(x_1, x_2)$ simply by sending x_2 to infinity, i.e.,

$$F(x_1) = \lim_{x_2 \rightarrow \infty} F(x_1, x_2). \quad (12)$$

In fact,

$$\lim_{x_2 \rightarrow \infty} F(x_1, x_2) = P(\{\omega : X_1(\omega) \leq x_1\} \cap \{\omega : X_2(\omega) \leq \infty\}) = P(\{\omega : X_1(\omega) \leq x_1\}) = F(x_1). \quad (13)$$

We can write the last equation in terms of PDFs as

$$\lim_{x_2 \rightarrow \infty} \int_{-\infty}^{x_1} \int_{-\infty}^{x_2} p(y_1, y_2) dy_1 dy_2 = \int_{-\infty}^{x_1} p(y_1) dy_1. \quad (14)$$

Since x_1 is arbitrary, it follows from (14) that

$$p(x_1) = \int_{-\infty}^{\infty} p(x_1, x_2) dx_2 \quad (\text{marginalization rule}). \quad (15)$$

Moreover, we have $F(\infty, \infty) = 1$, i.e.,

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p(x_1, x_2) dx_1 dx_2 = 1 \quad (\text{normalization condition}). \quad (16)$$

It is straightforward to extend these formulas to distribution functions and PDFs in more than two variables. For example, if $\mathbf{X}(\omega) = (X_1(\omega), X_2(\omega), X_3(\omega), X_4(\omega))$ is a four-dimensional random vector with distribution function $F(x_1, \dots, x_4)$ and PDF $p(x_1, \dots, x_4)$, then we can obtain the joint distribution function and the joint PDF of X_2 and X_3 , respectively, as

$$F(x_2, x_3) = F(\infty, x_2, x_3, \infty), \quad p(x_2, x_3) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p(x_1, x_2, x_3, x_4) dx_1 dx_4. \quad (17)$$

Example (Gaussian distribution): Consider the multivariate Gaussian PDF

$$p(x_1, \dots, x_n) = \frac{1}{\sqrt{(2\pi)^n \det(\Sigma)}} e^{-(\mathbf{x}-\boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}-\boldsymbol{\mu})/2}, \quad (18)$$

where

$$\mathbf{x}^T = [x_1 \ \dots \ x_n], \quad (19)$$

$$\boldsymbol{\mu}^T = [\mathbb{E}\{X_1\} \ \dots \ \mathbb{E}\{X_n\}] \quad (\text{mean}), \quad (20)$$

$$\Sigma_{ij} = \mathbb{E}\{X_i X_j\} - \mathbb{E}\{X_i\}\mathbb{E}\{X_j\} \quad (\text{covariance matrix}). \quad (21)$$

It is straightforward to show that all marginal PDF and distribution functions are still Gaussians of the form (18).

Independence. Let (Ω, \mathcal{B}, P) be a probability space. Two events $A \in \mathcal{B}$ and $B \in \mathcal{B}$ are said to be *independent* if the probability of their intersection (i.e., the probability that both events A and B happen) equals the product of their probabilities, i.e.,

$$A, B \in \mathcal{B} \text{ independent} \Leftrightarrow P(A \cap B) = P(A)P(B). \quad (22)$$

Consider now a random vector $\mathbf{X}(\omega) = (X_1(\omega), X_2(\omega))$ with components $X_1(\omega)$ and $X_2(\omega)$. We say that the random variables $X_1(\omega)$ and $X_2(\omega)$ are statistically independent if

$$P(\underbrace{\{\omega : X_1(\omega) \leq x_1\}}_{\text{event } A} \cap \underbrace{\{\omega : X_2(\omega) \leq x_2\}}_{\text{event } B}) = P(\{\omega : X_1(\omega) \leq x_1\})P(\{\omega : X_2(\omega) \leq x_2\}), \quad (23)$$

for all $x_1, x_2 \in \mathbb{R}$. Equation (23) can be written in terms of the cumulative distribution function as

$$F(x_1, x_2) = F(x_1)F(x_2). \quad (24)$$

This also implies that the joint PDF of X_1 and X_2 (if it exists) is simply the product of the PDF of X_1 and the PDF of X_2 , i.e.,

$$p(x_1, x_2) = p(x_1)p(x_2). \quad (25)$$

These formulas can be generalized to n independent random variables as

$$F(x_1, \dots, x_n) = F(x_1) \cdots F(x_n), \quad p(x_1, \dots, x_n) = p(x_1) \cdots p(x_n). \quad (26)$$

Examples:

- *Jointly uniform random vector.* Let \mathbf{X} be a n -dimensional random vector with zero-mean i.i.d. (independent identically distributed) uniform components in $[-1, 1]$. The joint PDF of \mathbf{X} is

$$p(x_1, \dots, x_n) = \begin{cases} \frac{1}{2^n} & (x_1, \dots, x_n) \in [-1, 1]^n \\ 0 & \text{otherwise} \end{cases} \quad (27)$$

- *Jointly normal random vector.* Let \mathbf{X} be a n -dimensional random vector with zero-mean i.i.d. Gaussian components with variance equal to one. The joint PDF of \mathbf{X} is

$$p(x_1, \dots, x_n) = \frac{1}{(2\pi)^{n/2}} e^{-\mathbf{x}^T \mathbf{x} / 2} \quad \mathbf{x} \in \mathbb{R}^n. \quad (28)$$

Clearly, from equation (18) we see that Gaussian random variables are independent if and only if

$$\mathbb{E}\{X_i X_j\} = \mathbb{E}\{X_i\}\mathbb{E}\{X_j\} \quad \text{for } i \neq j. \quad (29)$$

In general, if (29) is satisfied then we say that X_i and X_j are *uncorrelated*. Lack of correlation is a much weaker statement than independence, yet sufficient to claim independence for Gaussian random variables.

Expectation, joint moments, and joint cumulants. Let $\mathbf{X}(\omega) = (X_1(\omega), \dots, X_n(\omega))$ be a random vector defined on the probability space (Ω, \mathcal{B}, P) . For any measurable function $g(X_1, \dots, X_n)$ we define the expectation³ as

$$\mathbb{E}\{g(X_1, \dots, X_n)\} = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(x_1, \dots, x_n) p(x_1, \dots, x_n) dx_1 \cdots dx_n. \quad (31)$$

In particular, if $g(X_1, \dots, X_n) = X_1^{k_1} \cdots X_n^{k_n}$ then

$$\mathbb{E}\{X_1^{k_1} \cdots X_n^{k_n}\} = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} x_1^{k_1} \cdots x_n^{k_n} p(x_1, \dots, x_n) dx_1 \cdots dx_n \quad (\text{joint moments}) \quad (32)$$

The correlation matrix⁴ and the covariance matrix are defined as (see, e.g., (18))

$$\mathbb{E}\{X_i X_j\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_i x_j p(x_i, x_j) dx_i dx_j \quad (\text{correlation matrix}), \quad (34)$$

$$\mathbb{E}\{(X_i - \mu_i)(X_j - \mu_j)\} = \mathbb{E}\{X_i X_j\} - \mu_i \mu_j \quad (\text{covariance matrix}). \quad (35)$$

where $\mu_i = \mathbb{E}\{X_i\}$ (mean of X_i).

Remark: We say that two random variables $X_i(\omega)$ and $X_j(\omega)$ are *uncorrelated* if

$$\mathbb{E}\{X_i X_j\} = \mathbb{E}\{X_i\} \mathbb{E}\{X_j\}. \quad (36)$$

Independent random variables are always uncorrelated. In fact, let $p(x_i, x_j)$ be the joint PDF of X_i and X_j . We know that if X_i and X_j are independent then $p(x_i, x_j)$ can be factorized as

$$p(x_i, x_j) = p(x_i)p(x_j). \quad (37)$$

A substitution of (37) into (34) immediately yields (36). However, uncorrelated random variables are not necessarily independent.

We define the *moment generating function* of the random vector $\mathbf{X}(\omega) = (X_1(\omega), \dots, X_n(\omega))$ as

$$m(a_1, \dots, a_n) = \mathbb{E}\{e^{a_1 X_1 + \cdots + a_n X_n}\}. \quad (38)$$

³Note that the expectation $\mathbb{E}\{\cdot\}$ is a linear operator from a space of functions, e.g., the space of real-valued functions that are measurable with respect $p(x_1, \dots, x_n)$. Also, we do not need to assume the existence of the PDF to define the expectation operator. In fact, a more general expression for (31) is

$$\mathbb{E}\{g(X_1, \dots, X_n)\} = \int_{\Omega} g(X_1(\omega), \dots, X_n(\omega)) dP(\omega). \quad (30)$$

⁴Note that (34) follows from (32) using the marginalization property of the PDF. For instance

$$\begin{aligned} \mathbb{E}\{X_1 X_2\} &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} x_1 x_2 p(x_1, \dots, x_n) dx_1 \cdots dx_n \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 x_2 \left(\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} p(x_1, \dots, x_n) dx_3 \cdots dx_n \right) dx_1 dx_2 \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_1 x_2 p(x_1, x_2) dx_1 dx_2. \end{aligned} \quad (33)$$

It is straightforward to show that

$$\frac{\partial m(0, \dots, 0)}{\partial a_i} = \mathbb{E}\{X_i\}, \quad (39)$$

$$\frac{\partial^2 m(0, \dots, 0)}{\partial a_j \partial a_i} = \mathbb{E}\{X_i X_j\} \quad (40)$$

$$\begin{aligned} \frac{\partial^3 m(0, \dots, 0)}{\partial a_j \partial a_i \partial a_k} &= \mathbb{E}\{X_i X_j X_k\}, \\ &\dots \end{aligned} \quad (41)$$

Hence, the partial derivatives of the moment generating function evaluated at zero represent the joint moments of the components of random vector \mathbf{X} . Clearly, if $m(a_1, \dots, a_n)$ admits a convergent power series expansion at 0 then all joint moments exist.

The *cumulant generating function* of the random vector $\mathbf{X}(\omega) = (X_1(\omega), \dots, X_n(\omega))$ is defined as

$$\Psi(a_1, \dots, a_n) = \log(m((a_1, \dots, a_n))). \quad (42)$$

It is straightforward to show that

$$\frac{\partial \Psi(0, \dots, 0)}{\partial a_i} = \mathbb{E}\{X_i\}, \quad (43)$$

$$\frac{\partial^2 \Psi(0, \dots, 0)}{\partial a_j \partial a_i} = \mathbb{E}\{X_i X_j\} - \mathbb{E}\{X_i\}\mathbb{E}\{X_j\}, \quad (44)$$

$$\begin{aligned} \frac{\partial^3 \Psi(0, \dots, 0)}{\partial a_j \partial a_i \partial a_k} &= \mathbb{E}\{X_i X_j X_k\} - \mathbb{E}\{X_i\}\mathbb{E}\{X_j X_k\} - \mathbb{E}\{X_j\}\mathbb{E}\{X_i X_k\} - \mathbb{E}\{X_k\}\mathbb{E}\{X_i X_j\} \\ &\quad + 2\mathbb{E}\{X_i\}\mathbb{E}\{X_j\}\mathbb{E}\{X_k\}, \\ &\dots \end{aligned}$$

The quantities at the right hand side are known as *joint cumulants* of the random variables (X_1, \dots, X_n) . The cumulants are often denoted as $\langle X_i X_j \dots \rangle_c$ (see, e.g., [11])

$$\begin{aligned} \langle X_i \dots \rangle_c &= \mathbb{E}\{X_i\}, \\ \langle X_i X_j \dots \rangle_c &= \mathbb{E}\{X_i X_j\} - \mathbb{E}\{X_i\}\mathbb{E}\{X_j\}, \\ &\dots \end{aligned} \quad (45)$$

Characteristic function. The *characteristic function* of the random vector $\mathbf{X}(\omega) = (X_1(\omega), \dots, X_n(\omega))$ is defined as

$$\phi(a_1, \dots, a_n) = \mathbb{E} \left\{ e^{i(a_1 X_1 + \dots + a_n X_n)} \right\}. \quad (46)$$

Note that the characteristic function is the *Fourier transform*⁵ of the joint probability density function $p(x_1, \dots, x_n)$ and therefore it essentially carries the same information. The joint moments of \mathbf{X} can be computed as

$$\mathbb{E} \left\{ X_1^{k_1} \dots X_n^{k_n} \right\} = \frac{1}{i^{k_1 + \dots + k_n}} \frac{\partial^{k_1 + \dots + k_n} \phi(0, \dots, 0)}{\partial^{k_1} a_1 \dots \partial^{k_n} a_n}. \quad (47)$$

It is interesting to notice that the marginalization operation we have seen for the PDF, e.g.,

$$p(x_1) = \int_{-\infty}^{\infty} p(x_1, x_2, \dots, x_n) dx_2 \dots dx_n \quad (48)$$

⁵The Fourier transform (46) is taken in the appropriate function space, e.g., $L^2(\mathbb{R}^n)$ or in the space of tempered distributions $\mathcal{S}(\mathbb{R}^n)$.

turns out to be simplified substantially in Fourier space. Indeed

$$\phi(a_1) = \phi(a_1, 0, \dots, 0) = \mathbb{E} \{ e^{ia_1 X_1 + i0X_2 \dots + i0X_n} \}. \quad (49)$$

By using well-known series expansion of the complex exponential, it is possible to show that (see, e.g., [11])

$$\phi(a_1, a_2, \dots, a_n) = \exp \left[\sum_{\nu_1, \dots, \nu_n=0}^{\infty} \langle X_1^{\nu_1} \dots X_n^{\nu_n} \rangle_c \prod_{k=1}^n \frac{(ia_k)^{\nu_k}}{\nu_k!} \right], \quad (50)$$

where the series at the exponent excludes the case $\nu_1 = \dots = \nu_n = 0$. For example,

$$\phi(a_1, a_2) = \phi(a_1)\phi(a_2) \exp \left[\sum_{j,k=1}^{\infty} \langle X_1^j X_2^k \rangle_c \frac{(ia_1)^j (ia_2)^k}{j!k!} \right]. \quad (51)$$

Clearly, if X_1 and X_2 are independent we have $\langle X_1^j X_2^k \rangle_c = 0$ for all i and j and therefore (51) reduces to

$$\phi(a_1, a_2) = \phi(a_1)\phi(a_2). \quad (52)$$

Clearly, this equation is the Fourier transform of the PDF $p(x_1, x_2) = p(x_1)p(x_2)$, and shows that if X_1 and X_2 are independent both the joint PDF and the joint characteristic function can be factorized as a product of one-dimensional functions.

Joint PDF of m functions of n random variables. Let $\mathbf{X}(\omega) = (X_1(\omega), \dots, X_n(\omega))$ be a random vector with joint probability density function $p(x_1, \dots, x_n)$. Define

$$\begin{cases} Y_1 = g_1(X_1, \dots, X_n) \\ \vdots \\ Y_m = g_m(X_1, \dots, X_n) \end{cases} \quad (53)$$

What is the joint probability density function of the random vector $\mathbf{Y} = (Y_1, \dots, Y_m)$? Note that m can be smaller, equal or larger than n . These cases need to be handled differently.

- If $n = m$ and $\{g_1, \dots, g_m\}$ are distinct functions we proceed as in Theorem 1 below.
- If $m < n$ and $\{g_1, \dots, g_m\}$ are distinct functions we can add $m - n$ equations to complement the system so that we have n independent equations in n variables:

$$\begin{cases} Y_1 = g_1(X_1, \dots, X_n) \\ \vdots \\ Y_m = g_m(X_1, \dots, X_n) \\ Y_{m+1} = X_{m+1} \\ \vdots \\ Y_n = X_n \end{cases} \quad (54)$$

Once the joint PDF of Y_1, \dots, Y_n is known (using Theorem 1 below) then we can marginalize it with respect to (y_{m+1}, \dots, y_n) to obtain $p(y_1, \dots, y_m)$ as

$$p(y_1, \dots, y_m) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} p(y_1, \dots, y_m, y_{m+1}, \dots, y_n) dy_{m+1} \dots dy_n. \quad (55)$$

- If we have more equations than variables (i.e. $m > n$) then the computation of the joint PDF of (Y_1, \dots, Y_m) is not as straightforward as above. Consider for example the mapping $Y_1(\omega) = X(\omega)$ and $Y_2(\omega) = X^2(\omega)$. Here we have two functions of the same random variable. Note also that $Y_2 = Y_1^2$. It can be shown that the joint PDF of $Y_1 = X$ and $Y_2 = X^2$ is

$$p(y_1, y_2) = p_X(y_1) \delta(y_2 - y_1^2), \quad (56)$$

where p_X is the PDF of X and $\delta(\cdot)$ is the Dirac delta function.

Theorem 1. Let $\mathbf{x}_k(\mathbf{y})$ ($k = 1, \dots, r$) be the zeros of the nonlinear system of equations $\mathbf{y} = \mathbf{g}(\mathbf{x})$ defined in (53) (for $n = m$) or in (54) (for $m < n$). The joint PDF of Y_1, \dots, Y_n is given by

$$p_{\mathbf{Y}}(\mathbf{y}) = \sum_{i=1}^r \frac{p_{\mathbf{X}}(\mathbf{x}_i(\mathbf{y}))}{|J(\mathbf{x}_i(\mathbf{y}))|}, \quad (57)$$

where J is the Jacobian determinant⁶ associated with the mapping $\mathbf{g}(\mathbf{x})$ evaluated at $\mathbf{x}_i(\mathbf{y})$ (assumed non-zero).

The proof of this theorem is provided in [13, Chapter 8].

Example 1: Consider the mapping

$$Y_1 = X_1^2, \quad Y_2 = X_1 + X_2. \quad (59)$$

Suppose we know the joint PDF of X_1 and X_2 . What's the joint PDF of Y_1 and Y_2 ? The following mapping from (X_1, X_2) to (Y_1, Y_2) can be inverted as

$$\begin{cases} y_1 = x_1^2 \\ y_2 = x_1 + x_2 \end{cases} \Rightarrow \begin{cases} x_1 = \pm \sqrt{y_1} \\ x_2 = y_2 \mp \sqrt{y_1} \end{cases}. \quad (60)$$

The Jacobian determinant of (60) is easily obtained as

$$J(x_1, x_2) = \det \begin{bmatrix} 2x_1 & 0 \\ 1 & 1 \end{bmatrix} = 2x_1. \quad (61)$$

Hence, by applying Theorem 1, we obtain the following joint PDF of Y_1 and Y_2 is

$$p_{\mathbf{Y}}(y_1, y_2) = \frac{1}{2\sqrt{y_1}} [p_{\mathbf{X}}(\sqrt{y_1}, y_2 - \sqrt{y_1}) + p_{\mathbf{X}}(-\sqrt{y_1}, y_2 + \sqrt{y_1})] \quad y_1 \geq 0. \quad (62)$$

Example 2: Consider the mapping

$$Y_1(\omega) = X_1 \quad Y_2(\omega) = 2 \sin(2X_1(\omega) + X_2(\omega)), \quad (63)$$

where X_1 and X_2 are independent Gaussians with zero mean and variance one. In Figure 1 we estimate the joint PDF of Y_1 and Y_2 using the frequency approach, i.e., formula (11), and the 2D kernel density estimation method discussed in [2].

⁶In (57) it is assumed that

$$J(\mathbf{x}_i(\mathbf{y})) = \det \left[\frac{\partial \mathbf{g}(\mathbf{x})}{\partial \mathbf{x}} \right]_{\mathbf{x}=\mathbf{x}_i(\mathbf{y})} \neq 0. \quad (58)$$

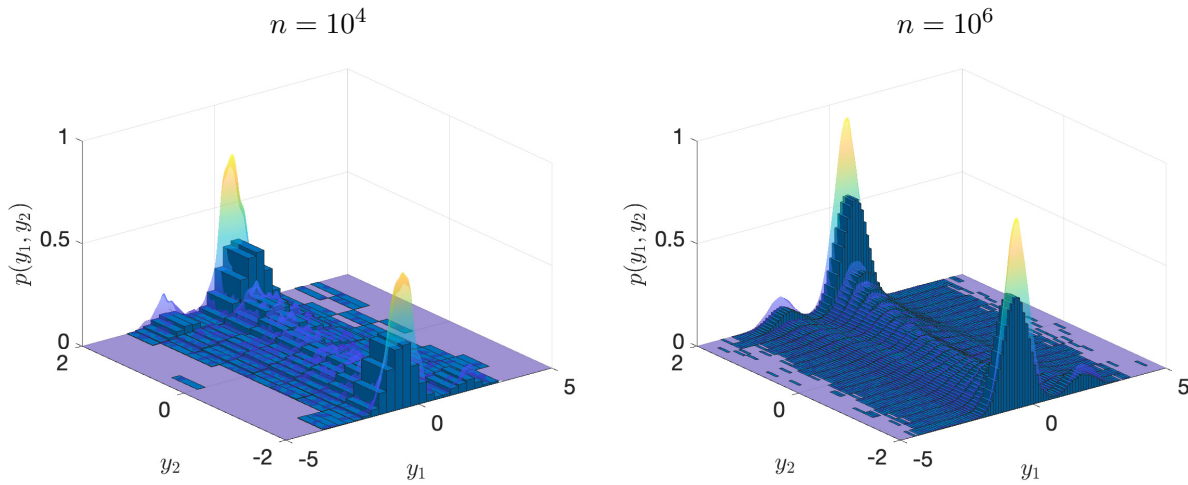


Figure 1: Estimation of the joint PDF of the random variables $Y_1 = X_1$ and $Y_2 = 2 \sin(2X_1 + X_2)$ where X_1 and X_2 are independent Gaussians with zero mean and variance one. We show the results we obtain with the frequency approach, i.e., formula (11) and the 2D kernel density estimation method discussed in [2] (transparent surface plot). We plot results for a different number of samples n .

Other methods to compute the joint PDF of functions of random vectors. There are of course other methods to compute the joint PDF (Y_1, \dots, Y_m) , given the joint PDF (X_1, \dots, X_n) . For instance, methods based on the Dirac delta function [8], or methods based on the joint characteristic function. With reference to the previous example we have the joint characteristic function

$$\phi_{\mathbf{Y}}(a_1, a_2) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{ia_1 x_1^2 + ia_2 (x_1 + x_2)} p(x_1, x_2) dx_1 dx_2. \quad (64)$$

Clearly, if $\phi_{\mathbf{Y}}(a_1, a_2)$ can be computed then we can simply inverse Fourier transform it to obtain the joint PDF of (Y_1, Y_2) . By using Dirac delta functions we can represent directly the joint PDF of the random variable

$$Y(\omega) = g(X_1(\omega), \dots, X_n(\omega)), \quad (65)$$

as

$$p(y) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \delta(y - g(x_1, \dots, x_n)) p(x_1, \dots, x_n) dx_1 \cdots dx_n \quad (66)$$

$$= \frac{1}{2\pi} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} e^{ia(y - g(x_1, \dots, x_n))} p(x_1, \dots, x_n) dx_1 \cdots dx_n da. \quad (67)$$

Example 3: Let $Y_1 = X$ and $Y_2 = X^2$ (two functions of one random variable). What is the joint PDF of Y_1 and Y_2 ? The mapping (53) yields a Jacobian determinant that is zero, and therefore the mapping is not invertible. This implies that theorem (1) cannot be applied. However, using the characteristic function approach we obtain

$$\phi(a_1, a_2) = \int_{-\infty}^{\infty} e^{ia_1 x + ia_2 x^2} p_X(x) dx. \quad (68)$$

Taking the inverse Fourier transform yields,

$$\begin{aligned}
 p(y_1, y_2) &= \frac{1}{(2\pi)^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{ia_1(x-y_1)+ia_2(x^2-y_2)} p_X(x) dx da_1 da_2 \\
 &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \delta(x-y_1) e^{ia_1(x^2-y_2)} p_X(x) dx da_2 \\
 &= \delta(y_1^2 - y_2) p_X(y_1).
 \end{aligned} \tag{69}$$

Example 4: If \mathbf{X} is Gaussian with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$, and \mathbf{A} is invertible then $\mathbf{Y} = \mathbf{A}\mathbf{X}$ is Gaussian with mean $\mathbf{A}\boldsymbol{\mu}$ and covariance $\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T$. To show this, we notice that we have the unique solution $\mathbf{X} = \mathbf{A}^{-1}\mathbf{Y}$. Hence by applying Theorem 1 we obtain the following PDF of \mathbf{Y}

$$\begin{aligned}
 p_{\mathbf{Y}}(\mathbf{y}) &= \frac{p_{\mathbf{X}}(\mathbf{A}^{-1}\mathbf{y})}{|\det(\mathbf{A})|} \\
 &= \frac{1}{\sqrt{2\pi}^n \det(\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)} \exp \left[-\frac{1}{2} (\mathbf{A}^{-1}\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{A}^{-1}\mathbf{y} - \boldsymbol{\mu}) \right] \\
 &= \frac{1}{\sqrt{2\pi}^n \det(\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)} \exp \left[-\frac{1}{2} (\mathbf{y} - \mathbf{A}\boldsymbol{\mu})^T (\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)^{-1} (\mathbf{y} - \mathbf{A}\boldsymbol{\mu}) \right].
 \end{aligned} \tag{70}$$

Remark: If (X_1, \dots, X_n) are independent random variables and (g_1, \dots, g_n) are n functions from \mathbb{R} into \mathbb{R} , then $Y_1 = g_1(X_1), \dots, Y_n = g_n(X_n)$ are independent random variables. It is straightforward to prove this statement using the Dirac delta function representation (or the characteristic function) of PDF mapping [8]. To this end, let

$$Y_i(\omega) = g_i(X_i(\omega)). \tag{71}$$

We have

$$\begin{aligned}
 p(y_1, \dots, y_n) &= \int_{-\infty}^{\infty} \prod_{j=1}^n \delta(y_j - g_j(x_j)) p(x_1, \dots, x_n) dx_1 \cdots dx_n \\
 &= \prod_{j=1}^n \int_{-\infty}^{\infty} \delta(y_j - g_j(x_j)) p(x_j) dx_j \\
 &= p(y_1) \cdots p(y_n).
 \end{aligned} \tag{72}$$

Sum of independent random variables. The PDF of the sum of independent random variables is the *convolution* the PDF of each variable. For example, let

$$Y = X_1 + X_2 + X_3 \tag{73}$$

be the sum of three independent random variables X_1 , X_2 and X_3 , with PDFs $p_1(x_1)$, $p_2(x_2)$ and $p_3(x_3)$ respectively. By using (66) we obtain

$$\begin{aligned}
 p(y) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \delta(y - x_1 - x_2 - x_3) p(x_1, x_2, x_3) dx_1 dx_2 dx_3 \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \delta(x_1 - y + x_2 + x_3) p_1(x_1) p_2(x_2) p_3(x_3) dx_1 dx_2 dx_3 \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p_1(x_2 + x_3 - y) p_2(x_2) p_3(x_3) dx_2 dx_3 \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p_1(x_1 - y) p_2(x_1 - x_3) p_3(x_3) dx_1 dx_3.
 \end{aligned} \tag{74}$$

In the last equality we considered the mapping $x_1 = x_2 + x_3$ as a coordinate change from x_1 to x_2 with parameter x_3 . Note that the process of computing the PDF of the sum of independent random variables can be also seen as a hierarchical process in which we proceed with two variables at a time. To this end, we first compute the PDF of $Z = X_2 + X_3$ as

$$p_Z(z) = \int_{-\infty}^{\infty} p_2(z - x_3)p_3(x_3)dx_3. \quad (75)$$

Clearly, Z is independent of X_1 and therefore the PDF of $Y = Z + X_1$ is

$$p_Y(y) = \int_{-\infty}^{\infty} p_1(y - x_1)p_Z(x_1)dx_1. \quad (76)$$

A substitution of (75) into (76) yields (74). A more direct proof relies on the characteristic function, which is the Fourier transform of (74). In fact

$$\begin{aligned} \phi_Y(a) &= \mathbb{E} \left\{ e^{ia(X_1+X_2+X_3)} \right\} \\ &= \mathbb{E} \left\{ e^{iaX_1} e^{iaX_2} e^{iaX_3} \right\} \\ &= \mathbb{E} \left\{ e^{iaX_1} \right\} \mathbb{E} \left\{ e^{iaX_2} \right\} \mathbb{E} \left\{ e^{iaX_3} \right\}. \end{aligned} \quad (77)$$

In the last step we used the fact that the expectation of a product of independent variables is the product of expectation. Equation (77) can be written as

$$\phi_Y(a) = \phi_{X_1}(a)\phi_{X_2}(a)\phi_{X_3}(a) \quad (78)$$

The inverse Fourier transform of (78) yields (74).

Example (sample average): Consider an experiment where we sample N independent realizations of a random variable $X(\omega)$ (like rolling a dice N times) and then take an average of all outcomes. Denote by $X_j(\omega)$ the outcome of the random variable $X(\omega)$ at the j -th sampling step and define

$$\bar{X}_N(\omega) = \frac{X_1(\omega) + X_2(\omega) + \cdots + X_N(\omega)}{N} \quad (\text{sample average}). \quad (79)$$

In this equation $X_j(\omega)$ are i.i.d. random variables with the same distribution as $X(\omega)$. Clearly if we repeat the sampling experiment multiple times, we have that $\bar{X}_N(\omega)$ attains different values (the samples of $\{X_1, \dots, X_N\}$ are different from experiment to experiment). As we will see when we study Monte Carlo sampling methods, $\bar{X}_N(\omega)$ is an approximation of $\mathbb{E}\{X\}$, i.e.,

$$\mathbb{E}\{X\} = \int_{-\infty}^{\infty} xp_X(x)dx \simeq \bar{X}_N(\omega) = \frac{X_1(\omega) + X_2(\omega) + \cdots + X_N(\omega)}{N}. \quad (80)$$

By using (77) and (78) it is straightforward to compute the characteristic function of $\bar{X}_N(\omega)$, and therefore its PDF. Specifically, we obtain⁷

$$\phi_{\bar{X}_N}(a) = \mathbb{E} \left\{ e^{iaX_j/N} \right\}^N \quad \text{for any fixed } j. \quad (81)$$

Let us provide two illustrative examples.

⁷Recall that the variables X_j are i.i.d. and therefore they have the same PDF and same characteristic function.

- *Sample average of a Gaussian random variable:* Consider a Gaussian random variable $X(\omega)$ with PDF

$$p_X(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}. \quad (82)$$

The sample average (79) in this case is a combination of i.i.d. variables X_j , all of which have PDFs (82). This implies that X_j/N are Gaussians with variance $1/N^2$, i.e.,

$$p_{X_j/N}(x) = \frac{N}{\sqrt{2\pi}} e^{-x^2 N^2/2}. \quad (83)$$

The characteristic function of X_j/N , i.e., $\mathbb{E}\{e^{iaX_j/N}\}$ is obtained by taking the Fourier transform of (83). This yields

$$\mathbb{E}\{e^{iaX_j/N}\} = e^{-a^2/(2N^2)} \quad (84)$$

Taking the product in (81) we obtain

$$\phi_{\bar{X}_N}(a) = e^{-a^2/(2N)} \quad \Leftrightarrow \quad p_{\bar{X}_N}(x) = \sqrt{\frac{N}{2\pi}} e^{-x^2/(2N)}. \quad (85)$$

Therefore, the sample average $\bar{X}_N(\omega)$ is distributed as a Gaussian random variable with variance $1/N$. This means that as we sum up more and more terms in (79) we have a concentration of measure phenomenon such that the Gaussian 85 gets more and more concentrated nearby 0. Moreover,

$$|\bar{X}_N(\omega) - \mathbb{E}\{X\}| = O\left(\frac{1}{\sqrt{N}}\right) \quad (86)$$

i.e., the sample mean $\bar{X}_N(\omega)$ converges⁸ to the mean of $X(\omega)$ at a rate $1/\sqrt{N}$.

- *Sample average of a Cauchy random variable:* Consider a Cauchy random variable $X(\omega)$ with PDF

$$p_X(x) = \frac{1}{\pi(1+x^2)} \quad (87)$$

We have seen that all moments of X (including the mean) are undefined, e.g.,

$$\int_{-\infty}^{\infty} xp_X(x)dx = \infty - \infty. \quad (88)$$

Does that mean that if we sample the Cauchy variable X and compute the sample mean (79) we may not converge to anything even for a very large number of samples N ? Yes! To show this, recall that the PDF of X_j/N in this case is

$$p_{X_j/N}(x) = \frac{N}{\pi(1+N^2x^2)}. \quad (89)$$

The Fourier transform of this PDF is

$$\phi_{X_j/N}(a) = \mathbb{E}\{e^{iaX_j/N}\} = e^{|a|/N}. \quad (90)$$

Therefore the characteristic function and PDF of the sample average (79) is

$$\phi_{\bar{X}_N}(a) = e^{|a|} \quad \Leftrightarrow \quad p_{\bar{X}_N}(x) = \frac{1}{\pi(1+x^2)}, \quad (91)$$

independently of the number of samples N ! Stated in different terms: no matter how many samples we consider, we have that the sample average of a Cauchy random variable is always a Cauchy random variable.

⁸As we will see, there are different modes of convergence of sequences of random variables, e.g., converge in probability, convergence in distribution, mean-square convergence, etc. In this case we have mean square convergence, which implies convergence in probability (thanks to the Markov inequality), and convergence in distribution.

Mapping random vectors to random vectors with desired distributions. It is always possible to flow a multivariate PDF $p_0(\mathbf{x})$ into another multivariate (target) PDF $p_1(\mathbf{x})$ using a time-dependent vector field, i.e., a system of ODEs of the form

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{f}(\mathbf{x}, t) \quad \mathbf{x}(0) = \mathbf{X}_0(\omega). \quad (92)$$

Indeed, as shown in Appendix A, the PDF of the state vector of the dynamical system (92) satisfies the Liouville equation

$$\frac{\partial p(\mathbf{x}, t)}{\partial t} + \nabla \cdot [\mathbf{f}(\mathbf{x}, t)p(\mathbf{x}, t)] = 0. \quad (93)$$

Of course, if we set $p(\mathbf{x}, 0) = p_0(\mathbf{x})$ and $p(\mathbf{x}, 1) = p_1(\mathbf{x})$ we can design (via optimization) a vector $\mathbf{f}(\mathbf{x}, t)$ that transports $p_0(\mathbf{x})$ to $p_1(\mathbf{x})$. This concept has been recently leveraged to develop samplers for arbitrary PDFs $p_1(\mathbf{x})$. The idea is to sample a known PDF $p_0(\mathbf{x})$, i.e., sample $\mathbf{X}_0(\omega)$ and then flow such samples to samples of $p_1(\mathbf{x})$ using (92) and an appropriate (optimized) vector field $\mathbf{f}(\mathbf{x}, t)$. Techniques such as continuous normalizing flows and *flow-matching* approaches [12, 6] are notable examples of methods in this category. Note that the target distribution can also be the distribution of a vector with independent components. In this sense, we can flow an arbitrary PDF into a fully separated one, hence transforming the components of an arbitrary random vector into independent random variables. In addition, the vector field $\mathbf{f}(\mathbf{x}, t)$ can be chosen to minimize a Wasserstein metric or other metrics (e.g. the KL divergence), in which case we talk about optimal mass transport [10].

Rosenblatt transformation. Another transformation that allows us to map a PDF $p_0(\mathbf{x})$ into another (arbitrary) PDF $p_1(\mathbf{x})$ is the *Rosenblatt transformation* [17]. Essentially, one can show that there exists a unique monotone increasing transformation of the form

$$\mathbf{T}(\mathbf{x}) = (T_1(x_1), T_2(x_1, x_2), \dots, T_n(x_1, x_2, \dots, x_n)) \quad (94)$$

such that

$$p_0(\mathbf{T}(\mathbf{x})) \det(\nabla \mathbf{T}(\mathbf{x})) = p_1(\mathbf{x}). \quad (95)$$

In this setting, the samples of $p_1(\mathbf{x})$ are obtained by simply mapping the samples of $p_0(\mathbf{x})$ via the transformation $\mathbf{T}(\mathbf{x})$ (when available). For the computation of $\mathbf{T}(\mathbf{x})$ see [17] and the references therein.

Mapping correlated Gaussian vectors to independent ones. To transform a correlated Gaussian vector into an independent one we just need a linear map. Essentially, if \mathbf{X} is Gaussian with zero mean and identity covariance then

$$\mathbf{Y} = \mathbf{A}\mathbf{X} + \boldsymbol{\mu} \quad (96)$$

has mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}^T$. This means that if we are interested in sampling a correlated Gaussian we can simply decompose the covariance $\boldsymbol{\Sigma}$ using the Cholesky factorization, sample \mathbf{X} and map the samples as in (96), where \mathbf{A} is the lower-triangular factor of the Cholesky decomposition. Similarly, if we to “de-correlate” a correlated Gaussian vector \mathbf{Y} we can apply invert the Cholesky factor of the covariance matrix and construct the vector

$$\mathbf{X} = \mathbf{A}^{-1}(\mathbf{Y} - \boldsymbol{\mu}). \quad (97)$$

Lebesgue spaces of random variables. The expectation operator $\mathbb{E}\{\cdot\}$ is a linear integral operator over a probability measure. Such an operator can be used to define norms and inner products in spaces of random variables. For example,

$$\mathbb{E}\{|X|^q\} = \int_{\Omega} |X(\omega)|^q dP(\omega) \quad q \in \mathbb{N} \quad (98)$$

is essentially a weighted q norm. The space of random variables satisfying $\mathbb{E}\{|X|^q\} < \infty$ is denoted as $L^q(\Omega, \mathcal{B}, P)$, in analogy with the classical Lebesgue space for functions. The case $q = 2$ is of particular importance as it has the structure of a Hilbert space. Specifically, for any two random variables in $L^2(\Omega, \mathcal{B}, P)$ we have the inner product

$$\mathbb{E}\{XY\} = \int_{\Omega} X(\omega)Y(\omega)dP(\omega) \quad (99)$$

and the norm

$$\mathbb{E}\{X^2\} = \int_{\Omega} X(\omega)^2 dP(\omega). \quad (100)$$

The inner product (99) allows us to define orthogonal random variables. Specifically, $X(\omega)$ and $Y(\omega)$ are orthogonal in $L^2(\Omega, \mathcal{B}, P)$ if they are uncorrelated, i.e., $\mathbb{E}\{XY\} = 0$. Also, $X(\omega)$ and $Y(\omega)$ are orthonormal if they are orthogonal and have norm equal to one, i.e., $\mathbb{E}\{X^2\} = \mathbb{E}\{Y^2\} = 1$.

Application to dynamical systems

Consider the following linear dynamical system

$$\begin{cases} \dot{x}(t) + \xi(\omega)x(t) = 0 \\ x(0) = x_0(\omega) \end{cases} \quad (101)$$

where $\xi(\omega)$ and $x_0(\omega)$ are independent random variables. Specifically $\xi(\omega)$ is uniformly distributed in $[0, 1]$, while $x_0(\omega)$ is Gaussian random variable with mean zero and variance one. As is well-known, the analytical solution of (101) is

$$x(t; \omega) = x_0(\omega)e^{-t\xi(\omega)}. \quad (102)$$

Let us compute the mean, the second-order moment and the auto-correlation function of the solution $x(t; \omega)$, i.e., $\mathbb{E}\{x(t; \omega)\}$, $\mathbb{E}\{x(t; \omega)^2\}$, and $\mathbb{E}\{x(t; \omega)x(t'; \omega)\}$ versus time. We have

$$\mathbb{E}\{x(t; \omega)\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x_0 e^{-x_0^2/2} dx_0 \int_0^1 e^{-t\xi} d\xi = 0, \quad (103)$$

$$\mathbb{E}\{x(t; \omega)^2\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x_0^2 e^{-x_0^2/2} dx_0 \int_0^1 e^{-2t\xi} d\xi = \frac{1}{2t} (1 - e^{-2t}), \quad (104)$$

$$\mathbb{E}\{x(t; \omega)x(t'; \omega)\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x_0^2 e^{-x_0^2/2} dx_0 \int_0^1 e^{-(t+t')\xi} d\xi = \frac{1}{t+t'} (1 - e^{-(t+t')}). \quad (105)$$

The one-time probability density function of $x(t; \omega)$ can be easily computed by using the Dirac delta function approach [8]. Indeed,

$$\begin{aligned} p(x, t) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \int_0^1 \delta(x - x_0 e^{-\xi t}) e^{-x_0^2/2} dx_0 d\xi \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \int_0^1 \frac{\delta(x_0 - x e^{\xi t})}{e^{-\xi t}} e^{-x_0^2/2} dx_0 d\xi \end{aligned} \quad (106)$$

$$= \frac{1}{\sqrt{2\pi}} \int_0^1 e^{\xi t - (x e^{\xi t})^2/2} d\xi. \quad (107)$$

Now consider the change of variables from ξ to u defined as

$$u = \frac{x e^{\xi t}}{\sqrt{2}} \Rightarrow d\xi = \frac{\sqrt{2}}{x t} e^{-\xi t} du. \quad (108)$$

A substitution of (108) into (107) yields

$$p(x, t) = \frac{1}{xt\sqrt{\pi}} \int_{x/\sqrt{2}}^{xe^t/\sqrt{2}} e^{-u^2} du \quad (109)$$

$$= \frac{1}{xt\sqrt{\pi}} \left[\operatorname{erf}\left(\frac{xe^t}{\sqrt{2}}\right) - \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right) \right]. \quad (110)$$

Liouville equation approach: We can transform the linear system (101) involving one random variable at the right hand side to an equivalent 2D linear system evolving from a random initial state (an no random variables at the right hand side). To this end, we notice that

$$\begin{cases} \dot{x}(t) + yx(t) = 0 \\ \dot{y}(t) = 0 \\ x(0) = x_0(\omega) \\ y(0) = \xi(\omega) \end{cases} \quad (111)$$

is completely equivalent to (101). In this setting, we can derive a linear transport equation for the joint PDF of $x(t; \omega)$ and $y(t; \omega)$, i.e., $x(t; \omega)$ and $\xi(\omega)$. Such PDF equation takes the form

$$\begin{cases} \frac{\partial p(x, y, t)}{\partial t} = \frac{\partial}{\partial x} (xyp(x, y, t)) + \frac{\partial}{\partial y} (xyp(x, y, t)) \\ p(x, y, 0) = p_{x_0}(x)p_{\xi}(y) \end{cases} \quad (112)$$

It can be verified by a direct substitution that the solution the initial value problem (112) is

$$p(x, y, t) = \frac{1}{\sqrt{2\pi}} e^{yt - (xe^{yt})^2/2}. \quad y \in [0, 1], \quad x \in \mathbb{R}. \quad (113)$$

Note that the joint PDF (113) was already obtained in equation (107), right before marginalizing with respect to ξ .

Conditional probability. Conditional probability is a measure of the probability of an event A occurring, given that another event B has already occurred. Such probability defined as⁹

$$P(A|B) = \frac{P(A \cap B)}{P(B)}. \quad (114)$$

Note that the conditional probability is non-zero if A and B are intersecting. Also note that if B is a subset of A then $P(A|B) = 1$.

Clearly, if A and B are independent events then by equation (22) we have that $P(A \cap B) = P(A)P(B)$. This implies that if A and B are independent then $P(A|B) = P(A)$. In other words, B has no effect whatsoever on the probability of A occurring. Moreover, $P(A \cap B) \leq P(B)$ (also $P(A \cap B) \leq P(A)$) and therefore we always have that

$$P(A|B) \leq 1. \quad (115)$$

⁹An example of conditional probability could be the following:

- Event A : “Daniele’s team scores a goal”.
- Event B : “Daniele takes a shot”.

The conditional probability $P(A|B)$, i.e., the probability that Daniele’s team scores a goal, conditional to Daniele taking a shot equals the probability that Daniele takes a shot and scores a goal, divided by the probability that Daniele takes a shot.

In the context of random vectors with multiple components, we may be interested in determining the conditional probability of an event involving one component, given that an event involving another component has already occurred. This yields the concept of *conditional CDF* and *conditional PDF*. Let us first clarify these concepts for a random vector with only two components $\mathbf{X}(\omega) = (X_1(\omega), X_2(\omega))$. By using the definition of the cumulative distribution function (5) we define (see [13, Ch. 7])

$$F(x_1|x_2) = \frac{F(x_1, x_2)}{F(x_2)} \Leftrightarrow F(x_1, x_2) = F(x_1|x_2)F(x_2). \quad (116)$$

Conditioning on particular outcomes. The determination of the conditional density of $X_1(\omega)$ given $X_2(\omega) = x_2$, i.e., when $X_2(\omega)$ attains a specific deterministic value x_2 , is of particular interest. This density cannot be derived directly from (114) or (116) because, as we know, the event $X_2(\omega) = x_2$ has zero probability (if X_2 is a continuous random variable). However, one can make sense of such conditional probability by taking a suitable limit. Specifically, consider

$$P(\underbrace{\{X_1(\omega) \leq x_1\}}_A \cap \underbrace{\{x_2 < X_2(\omega) \leq x_2 + \Delta x_2\}}_B) = F(x_1, x_2 + \Delta x_2) - F(x_1, x_2) \quad (117)$$

and

$$P(\{x_2 < X_2(\omega) \leq x_2 + \Delta x_2\}) = F(x_2 + \Delta x_2) - F(x_2). \quad (118)$$

In (117) it is understood that $F(x_1, x_2)$ is the joint distribution function of (X_1, X_2) , while in (118) $F(x_2)$ denotes the distribution function of X_2 alone. Clearly, for small Δx_2

$$F(x_1, x_2 + \Delta x_2) - F(x_1, x_2) \simeq \Delta x_2 \int_{-\infty}^{x_1} p(y_1, x_2) dy_1, \quad (119)$$

and

$$F(x_2 + \Delta x_2) - F(x_2) \simeq p(x_2) \Delta x_2. \quad (120)$$

Using (117)-(120) we can calculate the conditional probability

$$\frac{P(\{X_1(\omega) \leq x_1\} \cap \{x_2 < X_2(\omega) \leq x_2 + \Delta x_2\})}{P(\{x_2 < X_2(\omega) \leq x_2 + \Delta x_2\})} \simeq \frac{\Delta x_2 \int_{-\infty}^{x_1} p(y_1, x_2) dy_1}{\Delta x_2 p(x_2)}. \quad (121)$$

By sending Δx_2 to zero gives

$$F(x_1|X_2 = x_2) = \frac{\int_{-\infty}^{x_1} p(y_1, x_2) dy_1}{p(x_2)} \quad (\text{conditional CDF}). \quad (122)$$

Finally, by differentiating the previous equation with respect to x_1 we obtain

$$p(x_1|X_2 = x_2) = \frac{p(x_1, x_2)}{p(x_2)} \quad (\text{conditional PDF}). \quad (123)$$

In summary, to compute the conditional PDF, $p(x_1|X_2 = x_2)$ we literally take a section of the joint $p(x_1, x_2)$ for some fixed value of x_2 and then rescale the function we obtain by the number $p(x_2)$, i.e., the one-dimensional PDF of $p(x)$ of $X_2(\omega)$ evaluated at $x = x_2$. This procedure is illustrated in Figure 2 for a PDF represented in terms of a point cloud. Equation (123) can be written as

$$p(x_1, x_2) = p(x_1|x_2)p(x_2) = p(x_2|x_1)p(x_1) \quad (124)$$

which yields the identities

$$p(x_2) = \int_{-\infty}^{\infty} p(x_2|x_1)p(x_1)dx_1, \quad p(x_1) = \int_{-\infty}^{\infty} p(x_1|x_2)p(x_2)dx_2. \quad (125)$$

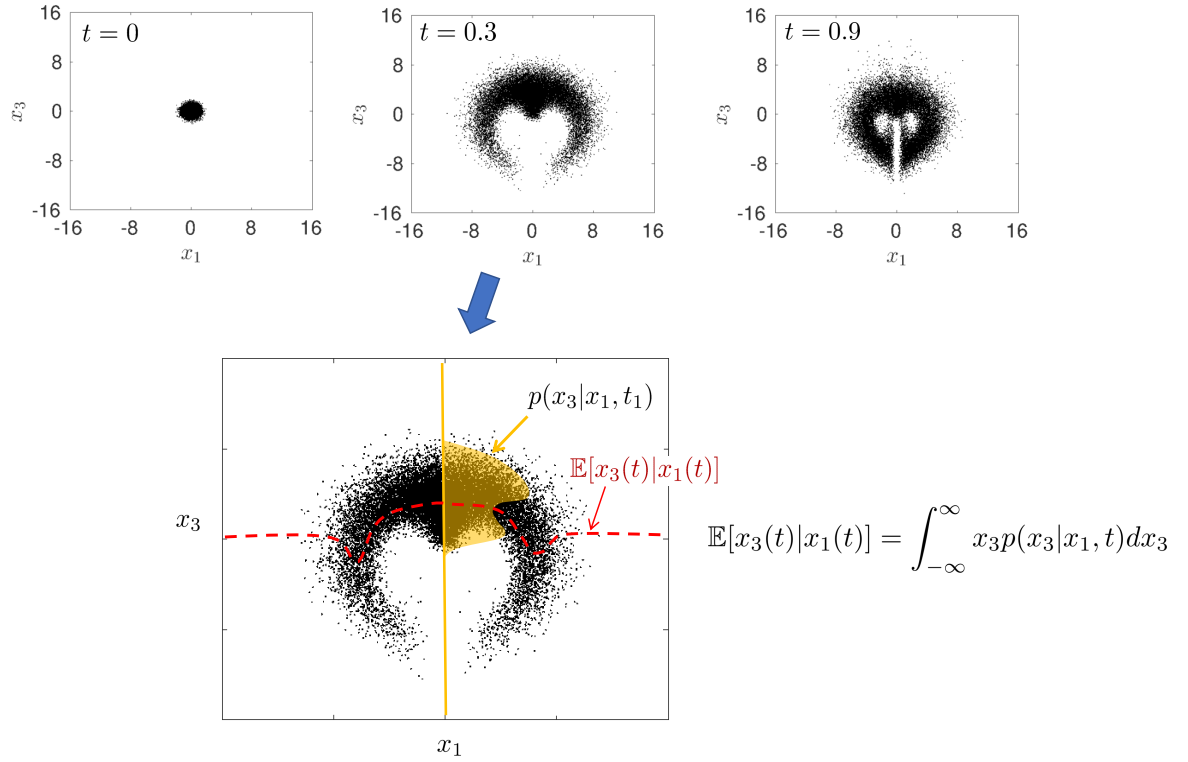


Figure 2: Point clouds representing the joint PDF of the phase variables $x_1(t)$ and $x_3(t)$ of the Kraichnan-Orzag system at different times, i.e., $p(x_3, x_1, t)$. Shown is the procedure to compute the conditional PDF $p(x_3|x_1, t)$ and the corresponding conditional mean $\mathbb{E}\{X_3|X_1 = x_1\}$.

The conditional probability density rule can be generalized to multiple random variables. For instance, if $p(x_1, x_2, x_3, x_4)$ denotes the joint PDF of four random variables then

$$p(x_1, x_2, x_3, x_4) = p(x_1|x_2, x_3, x_4)p(x_2, x_3, x_4) = p(x_1|x_2, x_3, x_4)p(x_2|x_3, x_4)p(x_3|x_4)p(x_4). \quad (126)$$

Moreover, conditional probability densities satisfy the *marginalization rule*. For instance

$$p(x_1, x_3|x_4, x_5) = \int_{-\infty}^{\infty} p(x_1, x_2, x_3|x_4, x_5) dx_2. \quad (127)$$

This property follows directly from the definition of conditional probability density (123).

Conditional expectation. Let $\mathbf{X}(\omega)$ and $\mathbf{Y}(\omega)$ be two random vectors defined on the probability space (Ω, \mathcal{B}, P) . The *conditional mean* of $g(\mathbf{X}(\omega))$ (g is a measurable function) assuming $\mathbf{Y}(\omega) = \mathbf{y}$ is defined as¹⁰

$$\mathbb{E}\{g(\mathbf{X})|\mathbf{Y} = \mathbf{y}\} = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(\mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x}, \quad (128)$$

where

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{y})} \quad (129)$$

is the conditional PDF of $\mathbf{X}(\omega)$ given $\mathbf{Y}(\omega) = \mathbf{y}$. Note that the $\mathbb{E}\{g(\mathbf{X})|\mathbf{Y} = \mathbf{y}\}$ is a function of \mathbf{y} . The conditional mean defined in equation (128) allows us to write the conditional moments of a random

¹⁰The conditional mean in equation (128) is often written as $\mathbb{E}\{g(\mathbf{X})|\mathbf{Y}\}$.

variable or a random vector, given information on another random vector. For example, the conditional mean and conditional correlation of \mathbf{X} given $\mathbf{Y}(\omega) = \mathbf{y}$ are defined as

$$\mathbb{E}\{X_i|\mathbf{Y} = \mathbf{y}\} = \int_{-\infty}^{\infty} x_i p(x_i|\mathbf{y}) dx_i, \quad (130)$$

$$\mathbb{E}\{X_i X_j|\mathbf{Y} = \mathbf{y}\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x_i x_j p(x_i, x_j|\mathbf{y}) dx_i dx_j. \quad (131)$$

The conditional mean of a system with two random variables is visualized in Figure 2. By combining (129), (128) and (31) we see that

$$\mathbb{E}\{g(\mathbf{X})\} = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \mathbb{E}\{g(\mathbf{X})|\mathbf{Y} = \mathbf{y}\} p(\mathbf{y}) d\mathbf{y}. \quad (132)$$

In this sense, $\mathbb{E}\{g(\mathbf{X})|\mathbf{Y} = \mathbf{y}\}$ can be interpreted as a random variable, i.e., a scalar function of the random variable \mathbf{Y} which, if averaged over $p(\mathbf{y})$, yields exactly $\mathbb{E}\{g(\mathbf{X})\}$.

Bayes' theorem. Let $p(\mathbf{x}, \mathbf{y})$ the joint probability density of two random vectors. Using conditional probabilities we have

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) \quad \Leftrightarrow \quad p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}. \quad (133)$$

This can be written, e.g., as

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{\int p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) d\mathbf{x}}. \quad (134)$$

Note that by using the marginalization rule we also have

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{y})p(\mathbf{y}) d\mathbf{y}. \quad (135)$$

Sampling high-dimensional PDFs using Markov Chain Monte Carlo (MCMC)

Markov Chain Monte Carlo (MCMC) refers to a class of methods that allow us to sample high-dimensional probability density functions [3]. In MCMC we construct a discrete Markov process that has a stationary PDF that coincides with the PDF of interest, i.e., the PDF we'd like to sample from. Hence, simulations of the Markov chain provide samples of the high-dimensional PDF we are interested in, once a transient called *burn-in phase* of the chain, is completed. To simulate the Markov chain it is common to use the Monte Carlo method, hence the name Markov Chain Monte Carlo. There are many different MCMC algorithms available to sample from high-dimensional PDFs. Perhaps, the simplest ones are

- Gibbs sampling (briefly described hereafter);
- Metropolis-Hastings algorithm (see [3]).

Gibbs sampling. Suppose we are given a three-dimensional PDF $p(x_1, x_2, x_3)$ and that the conditional PDFs $p(x_1|x_2, x_3)$, $p(x_2|x_1, x_3)$ and $p(x_3|x_1, x_2)$ are all available¹¹. To sample from $p(x_1, x_2, x_3)$ we proceed as follows:

¹¹Recall that to compute the conditional PDF $p(x_1|x_2, x_3)$ we literally set x_2 and x_3 in $p(x_1, x_2, x_3)$ equal to some number, say $x_2 = x_2^*$ and $x_3 = x_3^*$ and then normalize the one-dimensional function $p(x_1, x_2^*, x_3^*)$ so that the integral with respect to x_1 equals one.

1. Initialize $x_2 = x_2^{(i)}$ and $x_3 = x_3^{(i)}$. Here $x_2^{(i)}$ and $x_3^{(i)}$ are two real numbers. The superscript “ i ” is an integer number that labels the sample

$$\mathbf{X}^{(i)}(\omega) = \begin{bmatrix} x_1^{(i)}(\omega) & x_2^{(i)}(\omega) & x_3^{(i)}(\omega) \end{bmatrix} \quad i \in \mathbb{N}. \quad (136)$$

2. Sample a new $x_1^{(i+1)}$ from the one-dimensional conditional PDF $p(x_1|x_2^{(i)}, x_3^{(i)})$ (e.g., using the inverse CDF approach).
3. With the sample $x_1^{(i+1)}$ available, sample a new $x_2^{(i+1)}$ from the one-dimensional conditional PDF $p(x_2|x_1^{(i+1)}, x_3^{(i)})$.
4. With the sample $x_2^{(i+1)}$ available, sample a new $x_3^{(i+1)}$ from the one-dimensional conditional PDF $p(x_3|x_1^{(i+1)}, x_2^{(i+1)})$.
5. Update $x_j^{(i)} \leftarrow x_j^{(i+1)}$ for $j = 1, 2, 3$ and go back to point 2.

This algorithm allows us to compute $\mathbf{X}^{(i+1)}$ from $\mathbf{X}^{(i)}$ by sampling known one-dimensional conditional transition densities. To sample from such arbitrary one-dimensional transition densities we can use different methods. If the inverse cumulative distribution of each conditional PDF is known (or computable), then we are all set. In fact we can just sample a uniform PDF in $[0, 1]$ and then map such samples using the inverse cumulative distribution function. The mapping $\mathbf{X}^{(i)} \rightarrow \mathbf{X}^{(i+1)}$ defines a *random walk* in \mathbb{R}^3 . It can be shown that the stationary distribution of such random walk coincides with $p(x_1, x_2, x_3)$. In other words, after the *burn-in phase* is completed, i.e., for sufficiently large i , we have that $\mathbf{X}^{(i)}(\omega)$ are samples of the joint PDF $p(x_1, x_2, x_3)$.

Example (Gibbs’s sampling): Suppose we are interested in sampling the joint PDF

$$p(x_1, x_2) = K \sin^2(x_1 x_2), \quad (x_1, x_2) \in [0, 1]^2 \quad (137)$$

where

$$K = \frac{4}{2 - \text{Si}(2)} \quad \text{Si}(x) = \int_0^x \frac{\sin(t)}{t} dt \quad (138)$$

using Gibbs’ sampling. To this end, we first compute the marginals

$$p(x_1) = \frac{K}{2} \left(1 - \frac{\sin(2x_1)}{2x_1} \right), \quad (139)$$

$$p(x_2) = \frac{K}{2} \left(1 - \frac{\sin(2x_2)}{2x_2} \right), \quad (140)$$

and the conditionals

$$p(x_1|x_2) = \frac{p(x_1, x_2)}{p(x_2)} = \frac{4x_2 \sin^2(x_1 x_2)}{2x_2 - \sin(2x_2)}, \quad (141)$$

$$p(x_2|x_1) = \frac{p(x_1, x_2)}{p(x_1)} = \frac{4x_1 \sin^2(x_1 x_2)}{2x_1 - \sin(2x_1)}. \quad (142)$$

In Figure 3 we plot the PDF (137) and the samples we obtain from the Gibbs’ algorithm. To sample from (141)-(142) we use the 1D inverse CDF algorithm with numerically computed CDFs.

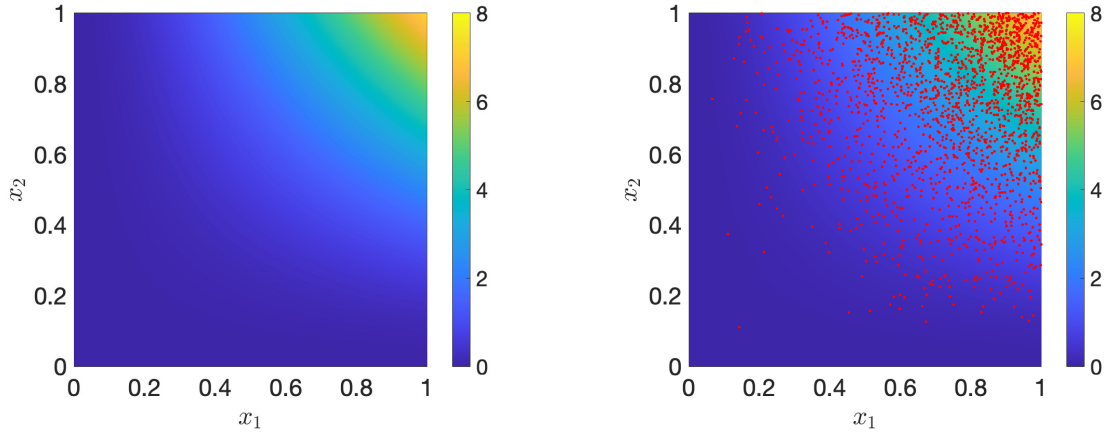


Figure 3: Plot of the 2D PDF (137) and samples of such PDF obtained by using Gibbs sampling.

Sampling approximate high-dimensional PDFs using copulas

A copula is a multivariate cumulative distribution function for which the marginal probability distribution of each variable is uniform on the interval $[0, 1]$. Copulas are used to model the dependence between random variables. Given a random vector

$$\mathbf{X}(\omega) = (X_1(\omega), \dots, X_n(\omega)) \quad (143)$$

we know can transform each component to a uniform random variable via the probability mappings

$$U_i(\omega) = F_i(X_i(\omega)) \quad i = 1, \dots, n. \quad (144)$$

A *copula* is defined as the cumulative distribution function of (U_1, \dots, U_n) , i.e.,

$$F_U(u_1, \dots, u_n) = P\{\{\omega : U_1(\omega) \leq u_1\} \cap \dots \cap \{\omega : U_n(\omega) \leq u_n\}\}. \quad (145)$$

The copula F_U is a CDF on the unit cube $[0, 1]^n$, and it contains all information on the dependence structure between the components of the random vector $\mathbf{X}(\omega)$. Note that by using the definition (144)

$$\begin{aligned} F_U(u_1, \dots, u_n) &= P\{\{\omega : U_1(\omega) \leq u_1\} \cap \dots \cap \{\omega : U_n(\omega) \leq u_n\}\} \\ &= P\{\{\omega : F_1(X_1(\omega)) \leq u_1\} \cap \dots \cap \{\omega : F_n(X_n(\omega)) \leq u_n\}\} \\ &= P\{\{\omega : X_1(\omega) \leq F_1^{-1}(u_1)\} \cap \dots \cap \{\omega : X_n(\omega) \leq F_n^{-1}(u_n)\}\} \\ &= F_{\mathbf{X}}(F_1^{-1}(u_1), \dots, F_n^{-1}(u_n)). \end{aligned} \quad (146)$$

Moreover, that equation (146) can be reversed to obtain

$$F_{\mathbf{X}}(x_1, \dots, x_n) = F_U(F_1(x_1), \dots, F_n(x_n)). \quad (147)$$

Hence, it appears that it is always possible write the joint CDF $F_{\mathbf{X}}(x_1, \dots, x_n)$ of a random vector \mathbf{X} in terms of a copula $F_U(u_1, \dots, u_n)$ and the individual CDFs F_i of each component $X_i(\omega)$. This result is known as *Sklar's theorem* (see, e.g., [5, Theorem 1.9]).

By differentiating (147) with respect to (x_1, \dots, x_n) yields the copula representation of the PDF of \mathbf{X} , i.e.,

$$p_{\mathbf{X}}(x_1, \dots, x_n) = p_U(F_1(x_1), \dots, F_n(x_n))p_1(x_1) \cdots p_n(x_n), \quad (148)$$

where $p_U(u_1, \dots, u_n)$ is the PDF of the copula.

Gaussian copula. The Gaussian copula is a distribution over the unit cuve $[0, 1]^n$. It is constructed from a multivariate normal distribution by using simple one-dimensional probability mappings. For a given correlation coefficient matrix \mathbf{R} (with entries in $[-1, 1]$) the Gaussian copula is defined as

$$p_{\mathbf{U}}(u_1, \dots, u_n) = \frac{1}{\sqrt{\det(\mathbf{R})}} \exp \left[\frac{1}{2} \mathbf{z}_{\mathbf{u}} (\mathbf{I} - \mathbf{R}^{-1}) \mathbf{z}_{\mathbf{u}}^T \right] \quad (149)$$

where

$$\mathbf{z}_{\mathbf{u}} = (F_g^{-1}(u_1), \dots, F_g^{-1}(u_n)) \quad (150)$$

and F_g denotes the Gaussian cumulative distribution function

$$F_g(u) = \frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{u}{\sqrt{2}} \right) \right] \quad (151)$$

Sampling from Gaussian copulas. Suppose we are interested in sampling a random vector with given marginals and correlation function

$$C_{ij} = \mathbb{E} \{ X_i X_j \} \quad (152)$$

In other words, we are not interested in sampling the full PDF of \mathbf{X} but rather create an approximate model, i.e., a *surrogate*, that allows us to sample the vector given only its marginals and the correlation matrix \mathbf{C} .

Gaussian copulas can be used for that as they have sufficient degrees of freedom to enforce a correlation structure while being consistent with marginals. The procedure to sample a Gaussian copula is very simple:

1. We first sample realizations of a Gaussian vector with zero mean and unit variance, say $\mathbf{Z}^{(j)}$.
2. Such sample vectors are then transformed to a vector with correlation \mathbf{R} (copula correlation) using the Cholesky factor¹² of \mathbf{R} . Call such samples $\mathbf{Y}^{(j)}$.
3. Each component of the rotated vectors is then mapped to a uniform distribution via the Gaussian probability mapping F_g (which is the same for all components), i.e., $\mathbf{U}^{(i)} = F_g(\mathbf{Y}^{(j)})$. The samples $\mathbf{U}^{(j)}$ are samples of the Gaussian copula.
4. At this point we generate a sample of the random vector of interest as \mathbf{X} by using the marginal CDF mapping as

$$\mathbf{X}^{(i)} = \left(F_1^{-1} \left(U_1^{(i)} \right), \dots, F_n^{-1} \left(U_n^{(i)} \right) \right). \quad (154)$$

Nataf transformation. It is important to emphasize that the correlation matrix \mathbf{R} is not the correlation matrix \mathbf{C} of the vector \mathbf{X} we are interested in sampling, but rather the “copula correlation matrix”. There is of course a relation between the two. Indeed, by using the the copula PDF we have

$$C_{ij} = \int_{[0,1]^2} F_i^{-1}(u_i) F_j^{-1}(u_j) p_{\mathbf{U}}(u_i, u_j; \rho_{ij}) du_i du_j \quad (155)$$

where

$$p_{\mathbf{U}}(u_i, u_j; \rho_{ij}) = \frac{1}{\sqrt{1 - \rho_{ij}^2}} \exp \left[\frac{1}{2(1 - \rho_{ij}^2)} \left[2\rho_{ij} F_g^{-1}(u_i) F_g^{-1}(u_j) - \rho_{ij}^2 (F_g^{-1}(u_i)^2 + F_g^{-1}(u_j)^2) \right] \right] \quad (156)$$

¹²The matrix \mathbf{R} is symmetric and positive definite. Hence it admits the Cholesky factorization $\mathbf{R} = \mathbf{L}\mathbf{L}^T$ where \mathbf{L} is lower triangular with positive diagonal entries. If \mathbf{Z} is a Gaussian vector with zero mean and unit covariance then $\mathbf{Y} = \mathbf{L}\mathbf{Z}$ is a Gaussian vector with zero mean and covariance \mathbf{R} . In fact

$$\mathbb{E} \{ \mathbf{Y}\mathbf{Y}^T \} = \mathbf{L} \mathbb{E} \{ \mathbf{Z}\mathbf{Z}^T \} \mathbf{L}^T = \mathbf{L}\mathbf{L}^T = \mathbf{R}. \quad (153)$$

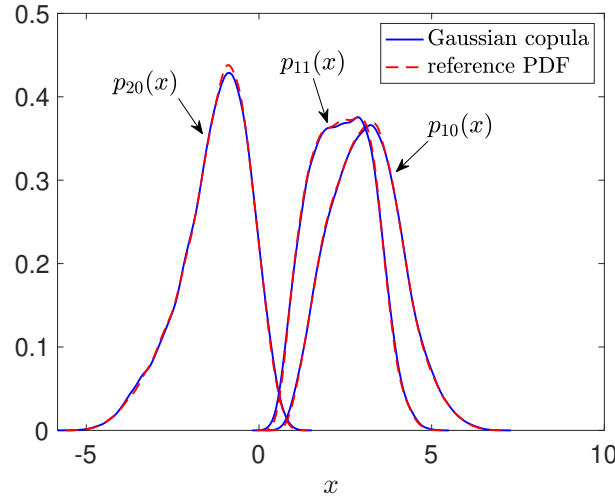


Figure 4: Marginals of the random function at points x_{10} , x_{11} and x_{20} . Comparison between reference marginals and marginals obtained from the Gaussian copula approximation. Note also that the marginals p_{10} and p_{11} overlap significantly due to the proximity of the spatial points x_{10} and x_{11} and the smoothness of the random function (159).

is the marginal of (149). Given a target correlation matrix C_{ij} , we can compute the copula correlation coefficients ρ_{ij} by solving the following sequence of 1D optimization problem

$$\min_{\rho_{ij} \in [-1, 1]} \left| C_{ij} - \int_{[0, 1]^2} F_i^{-1}(u_i) F_j^{-1}(u_j) p_U(u_i, u_j; \rho_{ij}) du_i du_j \right| \quad j > i. \quad (157)$$

Once all entries of \mathbf{R} , i.e., $\rho_{ij} \in [-1, 1]$ ($j > i$) are computed in this way, we may end up with a symmetric matrix \mathbf{R} that may not be positive definite. The closest (in the Frobenius norm) positive definite matrix to \mathbf{R} is easily obtained by doing a spectral decomposition and chopping off all negative eigenvalues.

$$\mathbf{R} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T \quad \Rightarrow \quad \mathbf{R}^+ = \mathbf{Q} \mathbf{\Lambda}^+ \mathbf{Q}^T \quad \mathbf{\Lambda}^+ = \max\{\mathbf{\Lambda}, \mathbf{0}\}. \quad (158)$$

The transformation between correlated Gaussian copula variable (with correlation \mathbf{R}) and the vector \mathbf{X} with given marginal distributions and correlation \mathbf{C} is called *Nataf transformation*.

Example (Gaussian copulas): Let us use Gaussian copulas to sample a random function $f(x; \omega)$ evaluated at $N = 100$ evenly-spaced grid points x_j . This yields a joint PDF in 100 dimensions. For this test problem, consider

$$f(x; \omega) = \sin(3x) e^{\xi_1(\omega) + \xi_2(\omega)} + e^{\cos(2x-2)} \xi_3(\omega) + \frac{2\xi_4(\omega)}{2\sin(2x) + 3} \quad x \in [0, 3\pi], \quad (159)$$

where (ξ_1, \dots, ξ_4) are four i.i.d uniform random variables in $[0, 1]$. The random vector representing $f(x; \omega)$ on the spatial grid has components

$$X_j = f(x_j; \omega), \quad x_j = \frac{3\pi(j-1)}{N-1}, \quad j = 1, \dots, 100. \quad (160)$$

We are given the marginals of $f(x; \omega)$, i.e., the PDFs of $f(x_i; \omega)$ at each spatial point x_i (see Figure 4), and the two-point correlation

$$C_{ij} = \mathbb{E} \{ f(x_i; \omega) f(x_j; \omega) \}. \quad (161)$$

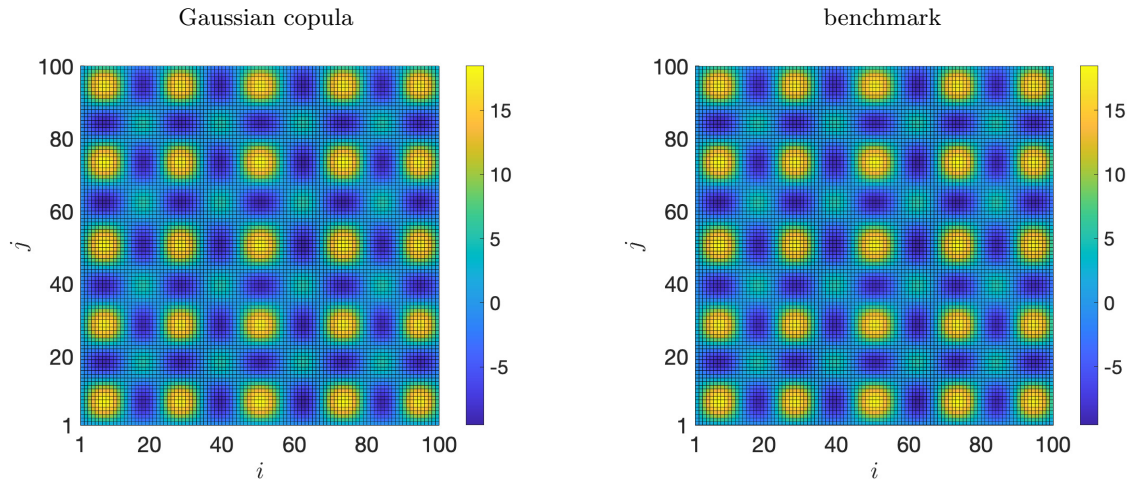


Figure 5: Comparison between the correlation matrix we obtain from Gaussian copula and the benchmark correlation (161).

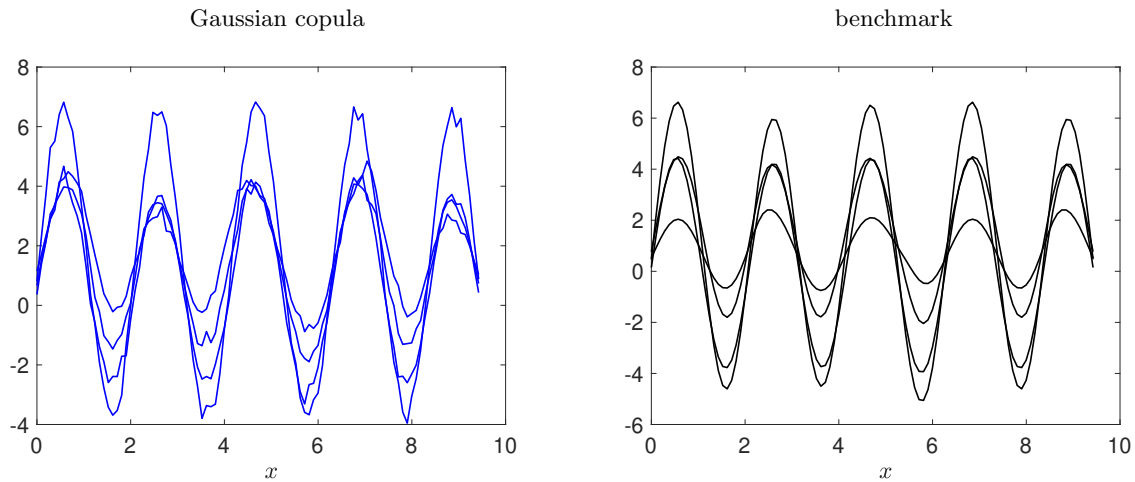


Figure 6: Samples of the Gaussian copula approximation (left) and samples from (159).

In Figure 5 we compare the correlation of the samples obtained from the Gaussian copula with the benchmark correlation. It is seen that the Gaussian copula has basically the same correlation and also the same marginals (see Figure 4) as the benchmark model.

In Figure 6 we plot a few samples of the Gaussian copula approximation and compare them with samples of the (159). Note that there exists a strong correlation structure in the random function (159) due to smoothness. Hence sampling each marginal distribution independently would completely destroy such correlation structure and result in sample patterns that are essentially Brownian motion (zig-zag) sort of noise. On the other hand, the sample patterns shown in Figure 6 capture the overall trend of the original patterns, with some minor deviations.

Learning high-dimensional PDFs from data

This is an active area of research today. There are two mainstreams commonly used for learning high-dimensional PDFs from data:

- Learning the functional form of the PDF using high-dimensional function representations such as tensors [18, 15], mixture PDF models, kernel density approaches, etc. These approaches can be

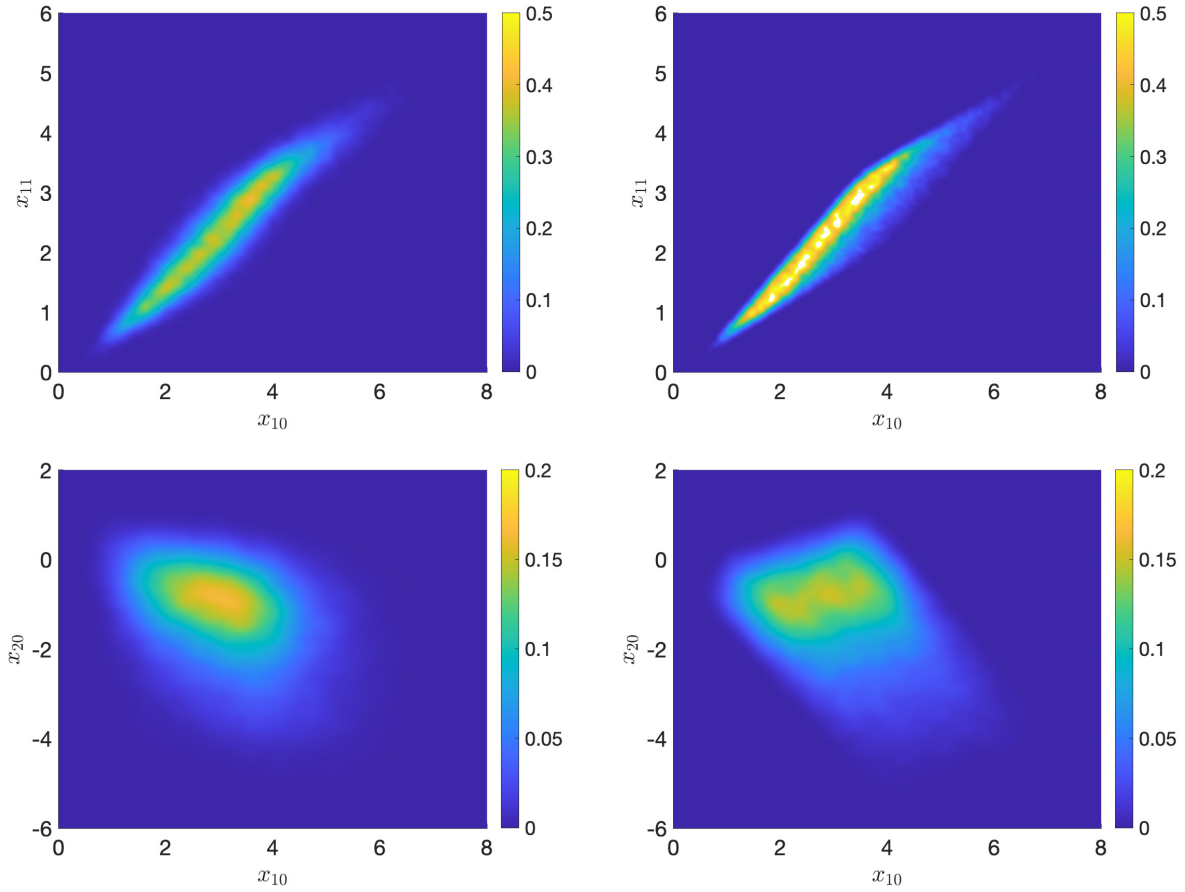


Figure 7: Joint PDFs from the Gaussian copula model (left), and benchmark PDFs (right). Note that the joint PDF of (159) at two neighboring nodes (x_{10} and x_{11}) is rather thin, due to the strong correlation between such random variables induced by the smoothness of $f(x; \omega)$. Note also that the Gaussian copula model does not reproduce exactly the joint PDF, but rather it approximates it.

computationally expensive due to the high-dimensionality of the function $p(x_1, \dots, x_n)$, which requires a representation.

- Developing a sampler for the PDF we are estimating from data, rather than attempting to compute the functional form of the PDF. These methods are computationally much more efficient than learning the functional form of the PDF. There are many different techniques that have been proposed/developed in the past decade or so to learn a sampler for a high-dimensional PDF that is learned from data.
 - a) **Generative Adversarial Networks (GANs)**. GANs consist of two components: a generator that produces fake data samples and a discriminator that tries to distinguish between real data and generated samples. The generator learns to create samples that mimic the real data distribution. GANs do not provide an explicit PDF but can generate samples from the high-dimensional distribution. They are widely used in tasks like image generation, video synthesis, and more (see [7, 1]).
 - b) **Variational Autoencoders (VAEs)**. VAEs combine neural networks with probabilistic models. The encoder maps input data into a latent space, and the decoder reconstructs the data from this latent space. The VAE is trained to maximize a lower bound on the likelihood of the data, thus learning an approximate distribution. The model does not produce an explicit PDF but provides a way to sample from the learned distribution.

- c) **Diffusion Models.** Diffusion models gradually convert data samples into noise through a diffusion process, then learn to reverse this process to generate new samples. Like GANs, they can implicitly learn high-dimensional distributions but do so through iterative denoising steps. In this class of models we find, e.g., DDPM [14], score-based models [16], and INDM [9].
- d) **Continuous normalizing flows (CNFs).** In CNFs, the goal is to transform a simple base distribution (like a Gaussian) into a complex target distribution through a sequence of invertible transformations, governed by differential equations (ODEs). A particular case of CNFs is flow matching [12], which simplifies this process by directly learning the optimal velocity field, which describes how data points flow from the base distribution to the target. Unlike traditional approaches that require learning complex dynamics over continuous time, flow matching trains the model to match the flow of data at discrete steps, making it more efficient. This method allows for the accurate estimation of high-dimensional PDFs, while also ensuring that the transformation remains invertible. Flow approaches are provably convergent [6, 4].

Appendix A: Derivation of the Liouville equation

Consider the nonlinear dynamical system

$$\begin{cases} \frac{d\mathbf{x}(t)}{dt} = \mathbf{f}(\mathbf{x}(t)) \\ \mathbf{x}(0) = \mathbf{x}_0(\omega) \end{cases} \quad (162)$$

where $\mathbf{x}_0(\omega)$ is a random vector with known joint probability density function $p_0(\mathbf{x})$. We know that if $\mathbf{f}(\mathbf{x})$ is continuously differentiable in \mathbf{x} then (162) admits a smooth flow $\mathbf{x}(t, \mathbf{x}_0(\omega))$, which is at least continuously differentiable in \mathbf{x}_0 . The flow is also continuously differentiable in t , i.e., $\mathbf{x}(t, \mathbf{x}_0(\omega))$ is a diffeomorphism in t . We are interested in determining an evolution equation for $p(\mathbf{x}, t)$, i.e., the probability density function of $\mathbf{x}(t, \mathbf{x}_0)$ at time t . To this end, consider the characteristic function representation of the PDF $p(\mathbf{x}, t)$

$$\phi(\mathbf{a}, t) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} e^{i\mathbf{a} \cdot \mathbf{x}(t; \mathbf{x}_0)} p(\mathbf{x}_0) d\mathbf{x}_0 = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} e^{i\mathbf{a} \cdot \mathbf{x}} p(\mathbf{x}, t) d\mathbf{x} \quad (163)$$

Differentiating with respect to t yields

$$\begin{aligned} \frac{\partial \phi(\mathbf{a}, t)}{\partial t} &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} i\mathbf{a} \cdot \frac{\partial \mathbf{x}(t, \mathbf{x}_0)}{\partial t} e^{i\mathbf{a} \cdot \mathbf{x}(t; \mathbf{x}_0)} p(\mathbf{x}_0) d\mathbf{x}_0 \\ &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} i\mathbf{a} \cdot \mathbf{f}(\mathbf{x}(t, \mathbf{x}_0)) e^{i\mathbf{a} \cdot \mathbf{x}(t; \mathbf{x}_0)} p(\mathbf{x}_0) d\mathbf{x}_0 \\ &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} i\mathbf{a} \cdot \mathbf{f}(\mathbf{x}) e^{i\mathbf{a} \cdot \mathbf{x}} p(\mathbf{x}, t) d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \frac{\partial}{\partial \mathbf{x}} (e^{i\mathbf{a} \cdot \mathbf{x}}) \cdot \mathbf{f}(\mathbf{x}) p(\mathbf{x}, t) d\mathbf{x} \\ &= - \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} e^{i\mathbf{a} \cdot \mathbf{x}} \nabla \cdot (\mathbf{f}(\mathbf{x}) p(\mathbf{x}, t)) d\mathbf{x}. \quad (\text{integrating by parts}) \end{aligned} \quad (164)$$

By using (163) and (164) we obtain

$$\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} e^{i\mathbf{a} \cdot \mathbf{x}} \left[\frac{\partial p(\mathbf{x}, t)}{\partial t} + \nabla \cdot (\mathbf{f}(\mathbf{x}) p(\mathbf{x}, t)) \right] d\mathbf{x} = 0, \quad \text{for all } \mathbf{a} \in \mathbb{R}^n, \quad (165)$$

which implies that the function between square bracket must be equal to zero for all \mathbf{x} and all t , i.e.,

$$\frac{\partial p(\mathbf{x}, t)}{\partial t} + \nabla \cdot (\mathbf{f}(\mathbf{x}) p(\mathbf{x}, t)) = 0 \quad (\text{Liouville equation}). \quad (166)$$

References

- [1] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein GAN. *arXiv:1701.07875*, pages 1–30, 2017.
- [2] Z. I. Botev, J. F. Grotowski, and D. P. Kroese. Kernel density estimation via diffusion. *Annals of Statistics*, 38(5):2916–2957, 2010.
- [3] S. Brooks, A. Gelman, G. Jones, and X.-L. Meng. *Handbook of Markov Chain Monte Carlo*. Chapman & Hall/CRC, 2011.
- [4] X. Cheng, J. Lu, Y. Tan, and Y. Xie. Convergence of flow-based generative models via proximal gradient descent in Wasserstein space. *arXiv:2310.17582*, 2024.
- [5] C. Czado. *Analyzing dependent data with vine copulas*. Springer, 2019.
- [6] Y. Gao, J. Huang, Y. Jiao, and S. Zheng. Convergence of continuous normalizing flows for learning probability distributions. *arXiv:2404.00551*, 2024.
- [7] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Su, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27:2672–2680, 2014.
- [8] A. I. Khuri. Applications of Dirac’s delta function in statistics. *Int. J. Math. Educ. Sci. Technol.*, 35(2):185–195, 2004.
- [9] D. Kim, B. Na, S. J. Kwon, D. Lee, W. Kang, and I.-C. Moon. Maximum likelihood training of implicit nonlinear diffusion models. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, 2024.
- [10] S. Kolouri, S. R. Park, M. Thorpe, D. Slepcev, and G. K. Rohde. Optimal mass transport: Signal processing and machine-learning applications. *IEEE Signal Processing Magazine*, 34(4):43–59, 2017.
- [11] R. Kubo. Generalized cumulant expansion method. *Journal of the Physical Society of Japan*, 17(7):1100–1120, 1962.
- [12] Y. Lipman, R. T. Q. Chen, H. Ben-Hamu, M. Nickel, and M. Le. Flow matching for generative modeling. *arXiv:2210.02747*, 2023.
- [13] A. Papoulis. *Probability, random variables and stochastic processes*. McGraw-Hill, third edition, 1991.
- [14] Y. Ren, H. Zhao, Y. Khoo, and L. Ying. High-dimensional density estimation with tensorizing flow. *Res. Math. Sci.*, 10:1–25, 2023.
- [15] Y. Ren, H. Zhao, Y. Khoo, and L. Ying. High-dimensional density estimation with tensorizing flow. *Res. Math. Sci.*, 10:1–25, 2023.
- [16] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. *arXiv:2011.13456*, pages 1–36, 2021.
- [17] A. Spantini, D. Bigoni, and Y. Marzouk. Inference via low-dimensional couplings. *Journal of Machine Learning Research*, 19:1–71, 2018.
- [18] X. Tang and L. Ying. Solving high-dimensional Fokker-Planck equation with functional hierarchical tensor. *Journal of Computational Physics*, 511:113110, 2024.