Quasi-random sampling for multivariate distributions via generative neural networks

Marius Hofert (with Avinash Prasad and Mu Zhu)

2022-04-20



Department of Statistics and Actuarial Science

1 Classical copula modeling

- General goal: Modeling a df H or X ~ H with continuous margins
 F₁,..., F_d (for possibly high-dimensional, computational applications).

 - Dynamic: X models joint innovations of an ARMA–GARCH process or increments of dependent geometric Brownian motions.
- Sklar's Theorem:
 - Analytical: $H(x) = C(F_1(x_1), \ldots, F_d(x_d))$ for a copula C, a distribution function with U(0, 1) margins.
 - ► Stochastic: Quantile transformation X = (F₁⁻¹(U₁),..., F_d⁻¹(U_d)) ~ H and probability transformation U = (F₁(X₁), ..., F_d(X_d)) ~ C.

Why: Useful from a computational point of view (e.g., estimation) or
 © Marius Hofert
 2

under asymmetric information (margins known, dependence unknown).

- Statistics: Instead of X, we now have $X = (X_1^{\top}, \dots, X_n^{\top})^{\top} \in \mathbb{R}^{n \times d}$.
 - 1) Given X, compute the *pseudo-observations* $U = (U_1^{\top}, \dots, U_n^{\top})^{\top} \in \mathbb{R}^{n \times d}$ with

$$U_{ij} = \hat{F}_{n,j}(X_{ij}) = \frac{1}{n+1} \sum_{k=1}^{n} \mathbb{1}_{\{X_{kj} \le X_{ij}\}} = \frac{R_{ij}}{n+1}$$

where $R_{ij} = \operatorname{rank}(X_{ij})$ among the component sample X_{1j}, \ldots, X_{nj} .



2) Based on U, one needs to fit, test and select an adequate copula C.

- 3) Simulate $U_1, \ldots, U_{n_{gen}}$ from the fitted C and estimate, for example:
 - $F_S^{-1}(u)$ where $S = X_1 + \dots + X_d = F_1^{-1}(U_1) + \dots + F_d^{-1}(U_d)$

Examples (QRM):

- Value-at-risk $\operatorname{VaR}_{\alpha}(S) = F_S^{-1}(\alpha)$ of the total loss S
- Expected shortfall $\mathrm{ES}_{\alpha}(S) = \mathbb{E}(S \mid S > F_S^{-1}(\alpha))$
- Expectations $\mu = \mathbb{E}(\Psi_0(\mathbf{X})) = \mathbb{E}(\Psi(\mathbf{U}))$ where $\Psi(\mathbf{u})$ is given by $\Psi_0(F_1^{-1}(u_1), \dots, F_d^{-1}(u_d))$ and so $\mu \approx \frac{1}{n_{\text{gen}}} \sum_{i=1}^{n_{\text{gen}}} \Psi(\mathbf{U}_i)$. Examples (MC):
 - With $\Psi(u) = \mathbbm{1}_{\{u_1 > u, ..., u_d > u\}}$ we obtain that

$$\mu = \mathbb{E}(\Psi(\boldsymbol{U})) = \mathbb{P}(U_1 > u, \dots, U_d > u),$$

so exceedance probabilities, e.g., the probability of a flood over multiple dikes or joint losses in a stock portfolio.

- With $\Psi(\boldsymbol{u}) = \max\{(F_1^{-1}(u_1) + \dots + F_d^{-1}(u_d))/d - K, 0\}$ we obtain that $\mu = \mathbb{E}(\Psi(\boldsymbol{U}))$ is the expected payoff of a European basket call option.

- Problems (with this classical modeling approach):
 - 1) Step 2) above: Finding an adequate copula model C for the pseudoobservations in more than two dimensions (oftentimes all are rejected)
 - 2) Large variance $\operatorname{var}(\frac{1}{n_{\text{gen}}}\sum_{i=1}^{n_{\text{gen}}}\Psi(U_i))$ of the Monte Carlo estimator in rare-event simulation scenarios.
- Ideas:
 - 1) Use (specific) neural networks (NNs) (\Rightarrow flexible dependence).
 - 2) Use quasi-random sampling for such NNs (\Rightarrow variance reduction).
- Outline:
 - ▶ We first present basics of quasi-random sampling under dependence.
 - We then study how NNs can be used as (pseudo-)random number generators (RNGs) from copulas.
 - ▶ We then address how NNs can be used as quasi-RNGs (QRNGs).
 - For simplicity, we mainly focus on the case where F₁ = ··· = F_d are U(0, 1) in what follows (except for an ES_α example).

2 Quasi-random numbers for copulas

- If $C = \Pi$ (independence), pseudo-random numbers (PRNs) can be replaced by quasi-random numbers (QRNs) for variance reduction.
- QRN sequences are low-discrepancy sequences $P_n = \{v_i\}_{i=1}^n$ (middle) with $D^*(P_n) = \sup_{z \in (0,1]^d} \left| \frac{\#\{i: v_i \in [0,z)\}}{n} - \lambda([0,z)) \right| \in O(n^{-1} \log^d n).$
- Example: $U_1, \ldots, U_{n_{gen}}$ (left), randomized Sobol' $V_1, \ldots, V_{n_{gen}}$ (right):



 The RQMC estimator ⁿ_{gen} ⁿ_{i=1} Ψ(V_i) is unbiased, fast and has a smaller variance (which can be estimated from repeated randomizations).
 Marius Hofert

- Question: How can we obtain QRNs from a general copula C?
- Idea: Could define $D_C^*(P_n) = \sup_{z \in (0,1]^d} \left| \frac{\#\{i: v_i \in [0,z)\}}{n} \mathbb{P}(U \in [0,z)) \right|$, but what's the order? And this does not lead to a construction principle. However, one can transform $V_1, \ldots, V_{n_{gen}}$ to samples from C; see Cambou et al. (2017).
- Inverse Rosenblatt transform \mathcal{R}^{-1} : Bijection to transform $U' \sim U(0,1)^d$ to $U \sim C$ (known as conditional distribution method (CDM) for sampling, a generalization of the inversion method to d > 1):

$$\begin{split} &U_1 = U_1', \\ &U_2 = C_{2|1}^{-1}(U_2' \mid U_1), \quad \text{(compare with } X_2 = F_2^{-1}(U_2') \sim F_2\text{)} \\ &U_j = C_{j|1,\dots,j-1}^{-1}(U_j' \mid U_1,\dots,U_{j-1}), \quad j \in \{2,\dots,d\}. \end{split}$$

Formula for the implementation (which needs to be inverted!):

$$C_{j|1,\dots,j-1}(u_j \mid u_1,\dots,u_{j-1}) = \frac{D_{j-1,\dots,1} C_{1,\dots,j}(u_1,\dots,u_j)}{D_{j-1,\dots,1} C_{1,\dots,j-1}(u_1,\dots,u_{j-1})}.$$

 \Rightarrow Inverse known analytically for Normal, t, Clayton (but can be slow). © Marius Hofert 7 • **Example:** 1000 PRNs (left) and QRNs (right) from a Clayton copula.



- Disclaimer: We (here) judge from such pictures but even if not visible, there may be variance reduction.
- Evaluating $C_{j|1,...,j-1}^{-1}$ is time-consuming even for normal, t, Clayton.
- For Gumbel copulas, C_{j|1,...,j-1} is tractable but C⁻¹_{j|1,...,j-1} is not.
- See Cambou et al. (2017) for a different approach, but conceptually (first component \Rightarrow frailty) and numerically challenging (α -stable qf).

3 PRNs for copulas via GMMNs

Idea: To address Problem 1) (flexible dependence), we consider

$$\mathbb{E}(\Psi(\boldsymbol{U})) \approx \frac{1}{n_{\text{gen}}} \sum_{i=1}^{n_{\text{gen}}} \Psi(\boldsymbol{U}_i), \tag{1}$$

$$\boldsymbol{U}_{i} = f_{\hat{\boldsymbol{\theta}}}(F_{\boldsymbol{Z}}^{-1}(\boldsymbol{U}_{i}')), \quad i = 1, \dots, n_{\text{gen}},$$
(2)

where

- $\blacktriangleright \ \boldsymbol{U}_1',\ldots,\boldsymbol{U}_{n_{\sf gen}}' \overset{\text{ind.}}{\sim} \mathrm{U}(0,1)^p \text{ (later: } \boldsymbol{V}_1,\ldots,\boldsymbol{V}_{n_{\sf gen}}\text{; randomized Sobol');}$
- $F_{Z}^{-1}(u) = (F_{Z_1}^{-1}(u_1), \dots, F_{Z_p}^{-1}(u_p))$ maps to a *prior distribution* (Φ);
- $f_{\hat{\theta}}$ is a trained generative neural network f_{θ} (GMMNs).
- Training of f_θ learns a map from pseudo-random numbers from Z (here: N_p(0, I_p)) to pseudo-random numbers from U ~ C.
- This is similar to R⁻¹, but computationally simpler to evaluate.
- Errors:
 - 1) Monte Carlo error (1) (can be made arbitrarily small by the SLLN)

2) Neural network (NN) "bottleneck" (2) (small if NN trained correctly) © Marius Hofert 9

3.1 What are NNs?



- Hyperparameters (fixed before training):
 - ► The \$\phi_l\$'s are the activation functions (e.g., ReLU \$\phi_l(x) = max{0, x}, sigmoid \$\phi_l(x) = 1/(1 + e^{-x})\$).
 - Number of layers, number of neurons per layer, number of epochs in training and batch size (more later).
- Parameter vector:

 $\boldsymbol{\theta} = (W_1, \dots, W_{L+1}, \boldsymbol{b}_1, \dots, \boldsymbol{b}_{L+1})$ consists of *weight matrices* and *biases* (initialized randomly and with zeros, respectively).

- θ is fitted (or: the NN is *trained*) by stochastic gradient descent
- The error is measured with a cost function E(U, Y), computed between
 - ▶ the training data $U = (U_1^{\top}, \dots, U_{n_{\text{trn}}}^{\top})^{\top}$ from C (= target) and
 - ► the outputs $Y = f_{\theta}(Z)$ of the NN from the prior sample $Z = (Z_1^{\top}, \dots, Z_{n_{\text{trn}}}^{\top})^{\top}$.
- NN applications often use a scaled MSE $E(U, Y) = \frac{1}{2n_{trn}} \sum_{i=1}^{n_{trn}} || \boldsymbol{u}_i \boldsymbol{y}_i ||_2^2$ as $E \Rightarrow$ Fails to learn the map from Z to U properly.

- Example:
 - ▶ Left: Subsample of size 2000 of 60 000 training data points of a t_{3.5} copula with correlation parameter 0.8.
 - Center: Sample of size 2000 of a NN trained with the MSE E.
 - ▶ **Right:** Sample of size 2000 of a NN trained with the "MMD" *E*.



 \Rightarrow The NN trained with the MSE E(U, Y) clearly did not capture the distribution correctly.

 \Rightarrow Learning a distribution (instead of a classification) is much harder.

3.2 What are GMMNs?

 Generative moment matching networks (GMMNs) use as E(U,Y) the sample maximum mean discrepancy MMD(U,Y)

$$\sqrt{\frac{1}{n_{\mathsf{trn}}^2}\sum_{i_1=1}^{n_{\mathsf{trn}}}\sum_{i_2=1}^{n_{\mathsf{trn}}}(K(\boldsymbol{U}_{i_1}, \boldsymbol{U}_{i_2}) - 2K(\boldsymbol{U}_{i_1}, \boldsymbol{Y}_{i_2}) + K(\boldsymbol{Y}_{i_1}, \boldsymbol{Y}_{i_2}))}$$

where K is a mixture of Gaussian kernels of different bandwidths. After experimentation, we chose

$$K(\boldsymbol{u}, \boldsymbol{y}) = \sum_{i=1}^{6} e^{-\frac{\|\boldsymbol{u}-\boldsymbol{y}\|_{2}^{2}}{2\sigma_{i}^{2}}} \quad \boldsymbol{\sigma} = (0.001, 0.01, 0.15, 0.25, 0.50, 0.75).$$

- Intuitively, MMD takes into account all pairs of observations between U_{i_1} and Y_{i_2} (desirable, but costly \Rightarrow mini-batch optimization).
- For the population version, one can show that MMD(U, Y) = 0 if and only if U = Y.

3.3 How are GMMNs trained?

- Training algorithm (mini-batch optimization):
 - 1) Initialize θ (weights W_1, W_2 : uniform entries; biases b_1, b_2 : 0)
 - 2) Partition the (n_{trn}, d) -data matrix U and prior (n_{trn}, d) -matrix Z into *batches* ($\approx n_{trn}/n_{bat}$ blocks of n_{bat} (batch size) rows each).
 - 3) For each batch, update θ by a stochastic gradient step (Adam).
 - 4) After $\approx n_{\rm trn}/n_{\rm bat}$ gradient steps, U and Z are exhausted and one epoch is completed. Shuffle their rows and go to Step 2).
 - 5) Training finishes after n_{epo} epochs.
- **Setup** in our experiments:
 - p = d (inspired from \mathcal{R}^{-1} in the CDM);
 - ▶ single hidden layer (universal approximation theorem: given suitable activation functions, single hidden layer NNs with d₁ < ∞ can approximate any continuous function on a compact subset of ℝ^d);
 - $d_1 = 300$ neurons in the hidden layer;

- ϕ_1 to be ReLU (fast) and ϕ_2 to be sigmoid (maps to (0,1));
- batch size n_{bat} = 5000 (trade-off: large enough for "bottleneck" error to be small; small enough to not be memory-prohibitive);
- number of epochs $n_{epo} = 300$; and (a rather large number of)
- $n_{\rm trn} = 60\,000$ pseudo-samples from C (available for all models).

3.4 How can GMMNs generate QRNs from C?

Algorithm 3.1 (GMMN quasi-random sampling)

- 1) Generate RQMC points $V_1, \ldots, V_{n_{gen}}$ (e.g., randomized Sobol').
- 2) Compute the prior samples $Z_i = F_Z^{-1}(V_i)$, $i = 1, ..., n_{gen}$.
- 3) Return the pseudo-observations of $Y_i = f_{\hat{\theta}}(Z_i)$, $i = 1, \dots, n_{gen}$.

GMMNs are fast to evaluate and sufficiently smooth to preserve low discrepancy. If all the mixed partial derivatives of $h = \Psi \circ f_{\hat{\theta}} \circ F_Z^{-1}$ exist a.e. and are continuous, then under Owen-type scrambling, $\operatorname{var}\left(\frac{1}{n_{\text{gen}}}\sum_{i=1}^{n_{\text{gen}}}h(\mathbf{V}_i)\right) = O(n_{\text{gen}}^{-3}(\log n_{\text{gen}})^{p-1}).$

3.5 Do GMMNs generate samples from C well?

Gumbel PRNs, GMMN PRNs, GMMN QRNs, *R*-transformed GMMN QRNs:



Gumbel-rotated t_4 mixture:



Marshall–Olkin copula $C(u_1, u_2) = \min\{u_1^{1-\alpha_1}u_2, u_1u_2^{1-\alpha_2}\}$ (without \mathcal{R}):



Nested Gumbel copula $C(u) = C_0(C_1(u_1, u_2), u_3)$ ($\tau \in \{0.25, 0.5\}$):



One can consider box plots of realization of the Cramér-von Mises statistic

$$S_{n_{\text{gen}}} = \int_{[0,1]^d} n_{\text{gen}} (C_{n_{\text{gen}}}(\boldsymbol{u}) - C(\boldsymbol{u}))^2 \, \mathrm{d}C_{n_{\text{gen}}}(\boldsymbol{u}),$$

where $C_{n_{\text{gen}}}$ is the *empirical copula* (ecdf of the pseudo-obs.) of $n_{\text{gen}} = 1000$ PRNs, GMMN PRNs and GMMN QRNs ($d \in \{3, 10\}$):



3.6 Variance-reduction effect?

Sample standard deviation of the estimator of $\mathbb{P}(U > 0.99)$ computed from PRNs, GMMN QRNs and QRNs:



Sample standard deviation of the estimator of $\mathbb{E}(S \mid S > F_S^{-1}(0.99))$ (N(0,1) margins) computed from PRNs and GMMN QRNs (no QRNs):



Summary

- We can learn a QRNG from any joint model based on a PRNG for C.
- If C is known, this is easy. If only a sample is available, n_{trn} needs to be sufficiently large (depending on the application).
- **Gain:** Universality (all models), computability (robustness, run time), especially useful for real data (where the true model is unknown).
- Challenges:
 - 1) Kendall's tau near 1;
 - 2) training needs a GPU server (evaluation "needs" TensorFlow);
 - 3) joint tail behavior;
 - 4) $d \gg 10$.
- **Open problem:** For $d \gg 10$, must training be improved (as distributions become harder to learn) or do RQMC point sets generally deteriorate? (GMMNs are still smooth). One could try (t, m, s)-nets other than Sobol' (but no related R package or standalone C code available).