

Conditional Loss and Deep Euler Scheme for Time Series Generation

Carl REMLINGER^{1,2,3}, Joseph MIKAEL^{2,3}, and Romuald ELIE¹

¹Université Gustave Eiffel

²EDF Lab

³FiME (Laboratoire de Finance des Marchés de l'Énergie)

Abstract

We introduce three new generative models for time series that are based on Euler discretization of Stochastic Differential Equations (SDEs) and Wasserstein metrics. Two of these methods rely on the adaptation of generative adversarial networks (GANs) to time series. The third algorithm, called Conditional Euler Generator (CEGEN), minimizes a dedicated distance between the transition probability distributions over all time steps. In the context of Itô processes, we provide theoretical guarantees that minimizing this criterion implies accurate estimations of the drift and volatility parameters. We demonstrate empirically that CEGEN outperforms state-of-the-art and GAN generators on both marginal and temporal dynamics metrics. Besides, it identifies accurate correlation structures in high dimension. When few data points are available, we verify the effectiveness of CEGEN, when combined with transfer learning methods on Monte Carlo simulations. Finally, we illustrate the robustness of our method on various real-world datasets.

1 Introduction

Time series Monte Carlo simulations are widely used for multiple industrial applications such as investment decisions (Kelliher and Mahoney, 2000), stochastic control (Pham, 2009) or weather forecasts (Mullen and Baumhelfner, 1994). They are notably considered in the financial sector, for market stress tests (Sorge, 2004), risk management and deep hedging (Buehler et al., 2019; Fecamp et al., 2020), or for measuring risk indicators such as Value at Risks (Jorion, 2000) among others. Providing Monte Carlo simulations representative of the time series of interest is a difficult and mostly manual task, which requires underlying modeling assumptions about the time dependence of the variables. Hence, it is difficult to update these models when a new type of data is observed, such as negative interest rates, negative electricity prices or unusual weather conditions. This naturally calls for the development of reliable model-free data generators for time series.

Generative methods such as Variational Auto Encoders (VAE) (Kingma and Welling, 2013) or Generative Adversarial Networks (GAN) (Goodfellow et al., 2014) provide state-of-the-art accuracy for the generation of realistic images (Xu et al., 2018) or text (Zhang et al., 2017). The development of similar generative methods for time series is very promising (Lyu et al., 2018; Chen et al., 2018). However, due to the complex and possibly non-stationary underlying temporal structure of the initial time series, these generative methods, especially GANs, are very difficult to apply as such (Yoon et al., 2019). Efficient generation of time series requires a proper learning of time-marginals as well as a faithfully representation of the underlying time structure.

In this paper, we embed time series into a discretized Euler approximations of Itô processes, which are characterized by their deterministic drift and volatility parameters. The three proposed generators rely on deep learning approximation of both drift and volatility functions. This representation benefits from a theoretically grounded temporal dynamic and provides a meaningful structure that avoids complex neural network architectures. Moreover, the considered Euler models allow tractable, at least controllable, generator outputs, which can be difficult with deep embedding such as (Yoon et al., 2019). This feature is a key

component in industrial applications, especially for decision-making-process.

By combining deep Euler representation with Wasserstein distance (Villani, 2008), we introduce the Euler Wasserstein GAN (EWGAN), inspired by (Arjovsky et al., 2017). Our second GAN-based-model, called Euler Dual Discriminator (EDGAN) is an adaptation of the DVDGAN presented in (Clark et al., 2019b). A spatial discriminator focuses on the accuracy of time-marginal distributions, while a temporal one focuses on the full sequence of generated time-series. Both methods compute the Wasserstein-1 distance and compete with the state-of-the-art algorithms Time Series GAN (TSGAN) (Yoon et al., 2019) and COTGAN (Xu et al., 2020) on both synthetic and real datasets. Nevertheless, all these GAN approaches still have difficulties to capture a proper temporal dynamics of the time series. We remedy to this problem by considering a loss function based on the conditional distributions $\mathcal{L}(X_{t+\Delta t} | X_t)$ of the generated time series. We introduce a Conditional Euler Generator (CEGEN) which optimizes a distance between the transition probability distributions at each time step. On the (large) class of Itô processes, we prove that minimizing this metric provides an accurate estimation of both the drift and volatility parameters.

A numerical study compares the three approaches to state-of-the-art GANs on synthetic and real datasets and shows the performances of our generators. We verify that our generators can learn to replicate Monte Carlo simulations of classical stochastic processes. Synthetic models give access to more reliable metrics (including theoretical), and allow to make connections between model-based Monte Carlo and model-free methods. EWGAN and EDGAN show a similar accuracy than TSGAN or COTGAN and capture more efficiently the time structure dynamics in dimension up to 20. The best performing model, CEGEN, is able to recover the underlying correlation (or independence) structure of time series, even in high dimensions. Moreover, we highlight the robustness of CEGEN, when combined with a transfer learning procedure when too few data are available. By properly mixing Monte Carlo generated and sparse real data during training, we CEGEN can take advantage of the synthetic simulations to improve its accuracy on generated samples.

Main Contributions:

- A theoretically grounded time series generator CEGEN combining an Euler structure with a dedicated loss on conditional distributions is proposed.
- Relying on a similar Euler structure, we also introduce two alternative GAN-like time series generators inspired by (Arjovsky et al., 2017) and (Clark et al., 2019b). They exhibit close performance to the state-of-the-art TSGAN (Yoon et al., 2019) and COTGAN (Xu et al., 2020) on marginal metrics, but capture more accurately the dynamic structure of Ito-based time series.
- A thorough numerical study on synthetic and various real world datasets demonstrate the robustness of our generators. Euler models succeed in correctly learning the underlying drifts and volatility structures of synthetic and outperforms the other considered methods on real datasets (accurate correlation structure up to dimension 20...). A transfer learning application when sparse data is available is provided.

2 Related works

The bootstrap method proposed by (Efron, 1982) is one of the first purely data-driven attempt to generate time series. Data samples are simply taken randomly with replacement. The scope of this technique is limited as it does not generate additional synthetic data but is based on historical ones. On the opposite, model-free approaches such as GAN allow to learn empirical distribution from data and thus to generate new samples. However, initial GAN proposals focused on the generation of non temporally ordered outputs. GAN’s architecture improvement for the time series case is an intensive area of research. For example, WaveGAN (Donahue et al., 2018) uses the causal architecture of WaveNet (Oord et al., 2016) for unsupervised synthesis of raw-waveform audio. Alternatively, several works consider recurrent neural networks to generate data sequentially and keep memory of the previous time series states (Mogren, 2016; Esteban et al., 2017).

Time Series GAN (TSGAN) (Yoon et al., 2019) introduces a state-of-the-art method for time series generation which stands out by its specific learning process. At each time step, an embedding network projects time series samples onto a latent space on which a GAN operates. TSGAN manages to get the correct marginal distributions and temporal correlation on classical processes and is used as a baseline in this paper. This method lacks of theoretical foundations ensuring a reliable quality of generated samples. As

the usage of generating model-free method grows rapidly, their application to sensitive fields (e.g. finance) must be considered cautiously and requires theoretical and empirical guarantees on the behavior of these generators. For this purpose, an active line of research looks towards reliable embedding of time series, such as signature (Fermanian, 2019; Buehler et al., 2019) or Fourier representation (Steinerberger, 2018).

Most recent applications on video generation focus on specific GAN architectures to capture the spatial-temporal dynamics. For example, MoCoGAN (Tulyakov et al., 2018) and DVD GAN (Clark et al., 2019b) combines two discriminators, one for the temporal dynamic and another one on each static frame. Specialized generator structures have also been designed, TGAN (Saito et al., 2017) proposed to generate a dynamic latent space and VGAN (Vondrick et al., 2016) combines two generators, one for marginals and another one for temporal dependencies. Following the idea of applying optimal transport to GANs (Arjovsky et al., 2017; Genevay et al., 2018), COTGAN (Xu et al., 2020) uses causal optimal transport for video sequence generation. To do so, the discriminant penalizes not-causal optimal transport plans, ensuring that the generator minimizes an adapted Wasserstein distance for time series. This approach benefits of solid theoretical foundations but still lacks of reliable empirical success for noisy time series generation.

3 Problem formulation

We aspire to design a time series generator which combines accurate estimation of time-marginal distributions while properly capturing temporal dynamics. The generator we propose is designed to be simple enough to be tractable (in the sense that outputs could be controlled) and theoretically grounded. To do so, we feed our algorithms with training time series data and seek to learn an empirical probability distribution that best approximates the data one. This task can be tricky, depending on the sequences lengths, the dimension, and the shape of the data distribution.

Although the idea of a model-free approach is attractive, we restrict ourselves to the context of Itô processes. This class of processes encompasses a wide range of time series and yet allows us to develop tractable models based on theory. In addition to providing a robust theoretical framework and controls on the processes generation, Itô processes allow to measure the accuracy of our generators on synthetic samples via closed form expressions or Monte Carlo simulators. In comparison to the classical literature (Wiese et al., 2020; Buehler et al., 2020), we do not assume the time series X to be stationary and allow ourselves to consider not-stationary sequences.

Itô process We are given i.i.d. samples of a time series, considered as a random vector $X = (X_{t_i})_{i=1\dots N}$ on $\mathbb{R}^{d \times N}$, starting from a point $X_0 \in \mathbb{R}^d$ and observed on a time grid $\mathcal{T} := \{0 = t_0 < t_1 < \dots < t_N = T\}$. For the sake of simplicity, in the following, we assume a regular time grid with mesh size Δt . The discrete time samples are supposed to be drawn from a continuous time underlying process X having the following Itô dynamics:

$$dX_t = b_X(t, X_t)dt + \sigma_X(t, X_t)dW_t, \quad (1)$$

where $b_X : \mathbb{R} \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the drift term, $\sigma_X : \mathbb{R} \times \mathbb{R}^d \rightarrow \mathcal{M}_{d \times d}$ the volatility term and W is a d -dimensional Brownian motion. The parameters b_X and σ_X are supposed to satisfy the usual Lipschitz conditions (Ikeda and Watanabe, 2014) ensuring existence and uniqueness of the solution of Eq.(1).

Deep Euler representation Samples of $X = (X_{t_i})_{i=1\dots N}$ are drawn from the continuous time Itô process with dynamics (1), and can approximately be viewed as samples drawn from the Euler discretization scheme of (1) given by

$$X_{t_i+\Delta t} = X_{t_i} + b_X(t_i, X_{t_i})\Delta t + \sigma_X(t_i, X_{t_i})\Delta W_{t_i}, \quad (2)$$

where $(\Delta W_{t_i})_i$ is a collection of i.i.d. $\mathcal{N}(0, \Delta t I_d)$ random variables. We rely on this approximation and introduce the following deep Euler representation of the time series. Starting at $t_0 = 0$, from $Y_0^\theta = X_0$ we generate time series by the following scheme:

$$Y_{t_i+\Delta t}^\theta = Y_{t_i}^\theta + b_Y^\theta(t_i, Y_{t_i}^\theta)\Delta t + \sigma_Y^\theta(t_i, Y_{t_i}^\theta)Z_{t_i}, \quad (3)$$

where Z_{t_i} are $\mathcal{N}(0, \Delta t I_d)$ i.i.d. random variables and the functions b_Y^θ and σ_Y^θ are θ -parametrized functions approximated by neural networks. Our objective is to learn b_Y^θ and σ_Y^θ , so that the distribution of the processes Y^θ and X are close.

Evaluation During the learning phase, neither b_X nor σ_X are given as inputs to any of the proposed models. However, the proposed formulation allows to compare *a posteriori* b_X and σ_X , when they are known, to the estimated b_Y^θ and σ_Y^θ . This provides a reliable metric on the generation accuracy. Moreover, this setup provides a convenient way to control the drift b_Y^θ and volatility σ_Y^θ functions. This task is delicate with deep embedding proposals. As already mentioned, it is highly challenging for generators of temporally ordered data to create samples accurate on time-marginal distributions as well as temporal dynamic metrics. In order to remedy to this weakness, we introduce below an innovative loss function based on a distance between conditional distributions at each time step.

4 Euler Generators

Euler Generators proposed in this paper are composed of two main elements: a network generating the drift and volatility terms of an Itô process and a distance between distributions to be minimized. The Itô structure facilitates the time series construction, while the distributional distance focuses on the law accuracy of the generated sequences. Both GAN-based and Deep Conditional methods described hereafter share this design.

4.1 Euler Generative Adversarial Networks

We propose two adaptations of GANs to time series that are based on the deep Euler representation presented in Eq.(3). The Wasserstein GAN (Arjovsky et al., 2017) seems to get rid of stability problems encountered in learning (mainly mode collapse) by adapting to the geometry of the underlying space. Our two GAN-based models are built on this proposal and both minimize (differentiable) Wasserstein-1 (\mathcal{W}_1) distance. The Rubinstein-Kantorovich duality allows to rewrite the \mathcal{W}_1 distance between two random variables Z_1 and Z_2 in the following way¹:

$$\mathcal{W}_1(\mathcal{L}(Z_1), \mathcal{L}(Z_2)) = \sup_{\|f\|_{\mathcal{L}} \leq 1} \mathbb{E}_{Z_1 \sim \mathcal{L}(Z_1)} [f(Z_1)] - \mathbb{E}_{Z_2 \sim \mathcal{L}(Z_2)} [f(Z_2)]. \quad (4)$$

Euler Wasserstein GAN (EWGAN) This model considers a Wasserstein GAN, where the generator relies on the Deep Euler representation (3) and optimizes the corresponding parameter θ . The discriminator d_φ parametrized by φ tries to find the optimal 1-Lipschitz function allowing to compute $\mathcal{W}_1(\mathcal{L}(X), \mathcal{L}(Y^\theta))$ using the Rubinstein-Kantorovich duality in (4). The 1-Lipschitz property of d_φ is guaranteed using the gradient penalty trick mentioned in (Gulrajani et al., 2017). The pseudocode of EWGAN is given in Alg.2 and details are provided in Appendix D. Overall, EWGAN minimizes the \mathcal{W}_1 distance between the distributions of the original $X = (X_t)_{t \in \mathcal{T}}$ and the generated one $Y^\theta = (Y_t^\theta)_{t \in \mathcal{T}}$:

$$\inf_{\theta} \mathcal{W}_1(\mathcal{L}(X), \mathcal{L}(Y^\theta)) = \inf_{\theta} \sup_{\varphi} \mathbb{E}_{X \sim \mathcal{L}(X)} [d_\varphi(X)] - \mathbb{E}_{Y^\theta \sim \mathcal{L}(Y^\theta)} [d_\varphi(Y^\theta)]. \quad (5)$$

Euler Dual Discriminator (EDGAN) Our second GAN-based-model, called EDGAN, is an adaptation of the Dual Video Discriminator (DVD) GAN (Clark et al., 2019a). DVD GAN uses attention networks and two discriminators in order to generate high fidelity videos. While the spatial discriminator focuses on time marginals and critics images in high resolution, the temporal one considers the full sequence of images in low resolution. We adapt these ideas to our context by considering in EDGAN, a similar dual discriminator architecture while the generator creates samples using the Deep Euler representation in (3). Our temporal discriminator follows a similar behavior as the one of EWGAN and focuses on $\mathcal{W}_1(\mathcal{L}(X), \mathcal{L}(Y^\theta))$. In the same time, our marginal discriminator focuses the computation of the W_1 distance between marginal distributions $\mathcal{W}_1(\mathcal{L}(X_t), \mathcal{L}(Y_t^\theta))$, for each $t \in \mathcal{T}$. The details and pseudocode of EWGAN are given in Alg.3 and provided in Appendix D.

¹ $\|f\|_{\mathcal{L}}$ denotes the smallest Lipschitz constant of f .

4.2 Conditional Loss Method

4.2.1 A loss function based on conditional distributions

The difficulty arising when trying to design a loss function for a time series generator comes from the need to get the correct balance between the marginal distribution fitness and the good representation of the temporal structure. On the one hand, we cannot only focus on marginals because having $\mathcal{L}(X_{t_i}) \sim \mathcal{L}(Y_{t_i}^\theta)$ for all $t_i \in \mathcal{T}$ does not imply that $b_X = b_{Y^\theta}$ nor that $\sigma_X = \sigma_{Y^\theta}$ (see the counterexample in Appendix A.1). On the other hand, instead of working on marginals, one can wonder if considering time series realization as samples of a vector defined on \mathbb{R}^{n+1} provides better results. Unfortunately, and as mentioned in (Yoon et al., 2019), learning the joint distribution $\mathcal{L}(X_{t_0}, \dots, X_{t_n})$ may not be sufficient to guarantee that the network captures the temporal dynamics, even with memory-based networks. An empirical example of unsatisfactory generation based on joint law is illustrated in Figure in Appendix 4, the generated trajectories are smooth. In the case of time series, one should simply refrain from applying a loss based only on marginal or joint distributions. To provide a reliable solution to this issue, we propose to focus on the transition probabilities at each time step by conditioning on the previous state. Moreover, by doing so, we are able to produce theoretical results on Itô coefficient estimation accuracy.

4.2.2 CEGEN Algorithm

Contrarily to the previous GAN-based generators, CEGEN does not require a discriminator network. The idea consists in considering a loss function that compares the conditional distributions $\mathcal{L}(Y_{t_{i+1}}^\theta | Y_{t_i}^\theta)$ with $\mathcal{L}(X_{t_{i+1}} | X_{t_i})$, for each time step $t_i \in \mathcal{T}$. The latter conditional distributions are Gaussian when considering Euler-discretized Itô processes. We consider the following metric:

$$\mathcal{W}_2^2(\mathcal{L}(X), \mathcal{L}(Y)) = \|\mathbb{E}[X] - \mathbb{E}[Y]\|_2^2 + \mathcal{B}^2(\text{Var}(X), \text{Var}(Y)) \quad (6)$$

where \mathcal{B} is the Bures metrics (Bhatia et al., 2019; Malago et al., 2018) defined by $\mathcal{B}^2(A, B) = \text{Tr}(A) + \text{Tr}(B) - 2\text{Tr}(A^{\frac{1}{2}}BA^{\frac{1}{2}})^{1/2}$, for positive definite matrices A and B . If X and Y are gaussian, \mathcal{W}_2 is the definition of the Wasserstein-2 distance (Gelbrich, 1990). This metric (6) captures meaningful geometric features between distributions, and \mathcal{W}_2 transportation plan is very sensitive to the outliers thus increases the distribution estimation accuracy. The Bures formulation (6) allows us to compute exactly the Wasserstein-2 distance, instead of regularized ones (Genevay et al., 2018; Cuturi, 2013).

Moreover, we want to have a theoretically grounded methodology and the Bures metric allows us to provide guarantees that minimizing the conditional loss implies accurate estimation for the drift and diffusion coefficients. We see that whenever the conditional distributions of the form $\mathcal{L}(X_{t_{i+1}} | X_{t_i} = z)$ and $\mathcal{L}(Y_{t_{i+1}} | Y_{t_i} = z)$ coincide in \mathcal{W}_2 , the drift and diffusion parameters coincide as well (see Prop. A.2 in Appendix). This is encouraging but in general conditioning from the very same point is complicated. Proposition 4.1 extends this property when the previous states belong to a small ball around z .

To build up our generator, we create at each time t_i a partition $(I_k)_{k \leq N_k}$ of the union of supports of X_{t_i} and $Y_{t_i}^\theta$. For a given batch of samples, $\mathcal{L}(X_{t_{i+1}} | X_{t_i} \in I_k)$ is approximated by extracting the elements $X_{t_{i+1}}$ such that $X_{t_i} \in I_k$. $\mathcal{L}(Y_{t_{i+1}} | Y_{t_i} \in I_k)$ is approximated in the same way. The \mathcal{W}_2^2 metric between the two conditional distributions are then summed up over all K subdivisions and over all time steps:

$$l(X, Y^\theta) = \sum_{i=0}^{N-1} \sum_{k=1}^{N_k} \mathcal{W}_2^2(\mathcal{L}(X_{t_{i+1}} | X_{t_i} \in I_k), \mathcal{L}(Y_{t_{i+1}}^\theta | Y_{t_i}^\theta \in I_k))$$

The pseudocode of CEGEN is given in Alg.1 and details are provided in Appendix D. Observe that our framework boils down to computing \mathcal{W}_2 metric between empirical distributions. Bures metrics is computed using the Newton-Schulz method (Muzellec and Cuturi, 2018), which is a differentiable way to get covariance matrix square roots.

4.2.3 Theoretical guarantee

In order to theoretically ground the choice of a loss function between conditional distributions based on \mathcal{W}_2 , we need to quantify how reducing the \mathcal{W}_2 loss (Eq. 6) implies proximity between drift and volatility

Algorithm 1 Algorithm CEGEN.

Input: \mathcal{D} samples of X , m batch size, K Nb of subdivisions, γ learning rate
Initialize: θ (randomly picked)
while Not converged **do**
 for $t_i = 0 \dots T$ **do**
 Sample m observations (x_{t_i+1}) from of X_{t_i+1}
 Sample $z \sim \mathcal{N}(0, I_D \Delta t)$
 $y_{t_i+1} \leftarrow y_{t_i} + g_\theta^b(t_i, y_{t_i}) \Delta t + g_\theta^\Sigma(t_i, y_{t_i}) z$
 $I_K \leftarrow K$ subdivisions of $\text{Supp}(X_{t_i}) \cup \text{Supp}(Y_{t_i})$
 for $k = 0 \dots K$ **do**
 $\ell_{t_i+1, k} \leftarrow \mathcal{W}_2^2(\mathcal{L}(x_{t_i+1} | x_{t_i} \in I_k), \mathcal{L}(y_{t_i+1} | y_{t_i} \in I_k))$
 end for
 end for
 $\theta = \theta - \gamma \nabla_\theta \sum_{t_i=1}^{T-1} \sum_{k=1}^K \ell_{k, t_i+1}$
end while
Output: y

parameters. An analysis on the topic is provided in Appendix A. The following result is allowed by the specific expression of the loss (6) implemented in CEGEN.

Proposition 4.1. *Assume that $\sigma_X^2(t_i, \cdot)$, $\sigma_{Y^\theta}^2(t_i, \cdot)$ are strictly positive and, together with $b_X(t_i, \cdot)$ and $b_{Y^\theta}(t_i, \cdot)$, K -Lipschitz in their second coordinate. For $t_i \in \mathcal{T}$, let $(I_k)_k$ be a regular partition covering $\text{Supp}(X_{t_i}) \cup \text{Supp}(Y_{t_i})$ with mesh size Δx . Let $\varepsilon > 0$.*

If $\mathcal{W}_2^2(\mathcal{L}(X_{t_i+1} | X_{t_i} \in I_k), \mathcal{L}(Y_{t_i+1}^\theta | Y_{t_i}^\theta \in I_k)) \leq \varepsilon^2$ for any k , then, for z in the partition

$$\|b_X(t_i, z) - b_{Y^\theta}(t_i, z)\|_2 \leq \frac{\varepsilon + \Delta x}{\Delta t} + 2K \Delta x.$$

Furthermore if $d = 1$,

$$\|\sigma_X(t_i, z) - \sigma_{Y^\theta}(t_i, z)\|_2 \leq \varepsilon / \sqrt{\Delta t} + 2K \Delta x.$$

and, when $d > 1$ and $\text{Tr}(\sigma_X^2(t_i, z)) = \text{Tr}(\sigma_{Y^\theta}^2(t_i, z)) = \alpha$, we have

$$\|\sigma_X(t_i, z) - \sigma_{Y^\theta}(t_i, z)\|_2 \leq \sqrt{\frac{2\alpha}{\Delta t}} \varepsilon + 2K \Delta x.$$

As described in A.3, the previous result is proved using useful inequalities between Hellinger and Bures distances. The α coefficient comes from the need of using density matrices, in practice one can easily normalize covariance matrices by their traces. Proposition 4.1 implies that by conditioning over sufficiently small intervals, a low \mathcal{W}_2 loss between transition distributions guarantees a good diffusion and drift representation.

5 Numerical Study

We now turn to the numerical evaluation of EWGAN, EDGAN and CEGEN in comparison to the state-of-the-art TSGAN and COTGAN, on various synthetic and real time series. Neural network architectures and hyper-parameters are described in Appendix D.

5.1 Datasets

Two kinds of datasets are used: synthetic and real time series dataset. In single dimension, we use Black&Scholes (BS) model ($dX_t = rX_t dt + \sigma X_t dW_t$) and an Ornstein-Uhlenbeck (OU) model ($dX_t = \theta(\mu - X_t) dt + \sigma dW_t$). For these two stochastic models, our empirical references are drawn from Monte Carlo (MC) simulators. The simulations are performed on a regular time grid of 30 dates, the maturity is 0.25 (1 simulation per day for 3 months) and $X_0 = 0.2$. BS model (resp. OU) has coefficients of $r = 0.8$, $\sigma = 0.3$ (resp. $\sigma = 0.1$, $\mu = 0.6$ and $\theta = 7$). In higher dimensions, we proceed with the same methodology but with multivariate correlated BS time series ($d = 4, 10, 20$). The real datasets include various nature of time series and are detailed in Appendix F.

5.2 Evaluation metrics

We consider several metrics to evaluate the accuracy of the generators. For all metrics the lower, the better.

(1) Marginal metrics. These metrics quantify the quality at each time step of the marginal distributions induced by the generated samples in comparison to the empirical ones. This includes Fréchet Inception Distance (FID) (Heusel et al., 2017) as well as classical statistics (mean, 95% and 5% percentiles, maximum and minimum denoted respectively Avg, q95, q05, Max, Min). We systematically compute the mean squared error (MSE) over time of these statistics between the real and generated samples. This helps measuring whether a generator manages to get an accurate overall envelope of the processes.

(2) Temporal dynamics. This metric aims at quantifying how the generator is able to capture the underlying time structure of the signal. For this purpose, we compute the difference between the quadratic variations of both reference and generated time series. The quadratic variation (QVar) of an Itô process X is given by $[X]_t = \int_0^t \sigma_X^2(s, X_s) ds$. Thus the temporal metric ensures that the diffusion σ_X is well estimated too. We compute $[X]_t$ in the discrete case with $\sum_i |X_{t_{i+1}} - X_{t_i}|^2$.

(3) Correlation structure. The metric denoted *Corr* in the following is the term-by-term MSE between empirical correlation from reference samples on one side and from generated samples on the other side. It evaluates the ability of a generator to capture the multi-dimensional structure of the signal.

(4) Underlying process parameters. A by-product output of Euler-based generators are the estimated drift $b_Y^\theta(\cdot)$ and diffusion $\sigma_Y^\theta(\cdot)$ coefficients of the generator. When using synthetic data, we can compare the true underlying processes parameters to the estimated ones. In the BS case, the drift and volatility coefficients are estimated by the empirical average of $(b_Y^\theta(t, Y_t^\theta)/Y_t^\theta)_{t \in \mathcal{T}}$ and $(\sigma_Y^\theta(t, Y_t^\theta)/Y_t^\theta)_{t \in \mathcal{T}}$. In the OU case, σ_Y^θ is estimated in a similar manner, while θ and μ are estimated by regressing b_Y^θ on (t, Y_t^θ) . These statistics cannot be computed in the same way with TSGAN due to its specific deep embedding, nor COTGAN.

(5) Discriminative and predictive scores. We use two distinct scores, as proposed in (Yoon et al., 2019). First, we train a classification model (a 2-layer LSTM) to distinguish real sequences from the generated ones. The accuracy of the classifier provides the discriminative score. Second, the predictive score is obtained by training a sequence-prediction model (a 2-layer LSTM) on generated time series to predict the next time step value over each input sequence. Performance is measured in terms of MAE.

5.3 One-dimensional simulated process (Exp. A)

We start with a unidimensional time series by comparing the five generators in the OU case. Figure 1 illustrates how crucial is the balance between the estimation of the marginal distributions and the temporal structure. On the one hand, the trend and marginal distributions of the time series generated by both GANs seem close to the empirical benchmark. However, the temporal dynamics between two time steps is not respected as confirmed by the QVar metric in Table 5 (in Appendix). On the other hand, CEGEN model manages to capture the overall envelope and is able to fit the dynamics of time series as the QVar metrics highlights in Table 5. Table 1 reports the reference drift and volatility coefficients with those obtained by the

	CEGEN	EWGAN	EDGAN
Black-Scholes			
\hat{r} (0.8)	0.739	0.581	0.996
$\hat{\sigma}$ (0.3)	0.324	0.314	0.379
Ornstein-Uhlenbeck			
θ (7.0)	7.05	4.36	4.68
$\hat{\mu}$ (0.6)	0.60	0.75	0.72
$\hat{\sigma}$ (0.1)	0.11	0.16	0.02

Table 1: **Exp. A.** Model parameter estimations for drift and volatility function.

three Euler-based generators. We can see a good estimation of the CEGEN method and to a lesser extent of the EWGAN method while EDGAN fails to estimate the parameters correctly. Euler structure alone does not manage to recover the right parameter values. To conclude this section, it appears that regarding the overall dynamics and the marginals, CEGEN is a reliable generator of time series. The question we address in the next section is how CEGEN scales to higher dimensions.

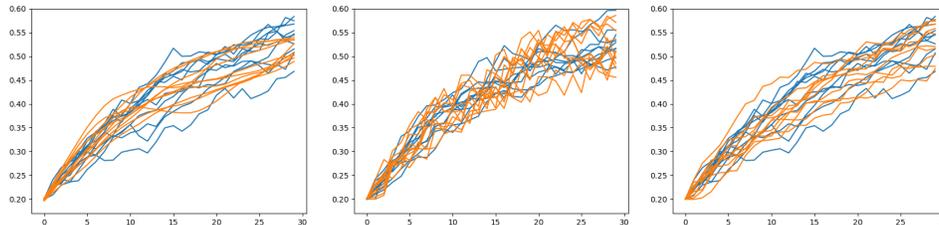


Figure 1: **Exp. A.** Example of Ornstein-Uhlenbeck samples (in blue) with COTGAN, TSGAN and CEGEN generations (orange).

5.4 Scaling the dimension (Exp. B)

Dim	CEGEN	EWGAN	EDGAN	TSGAN	COTGAN
4	.007	.015	.053	.177	.031
10	.011	.055	.022	.259	.035
20	.006	.034	.014	.481	.019

Table 2: **Exp. B** MSE between reference and generator empirical correlation matrices on Black-Scholes.

Table 2 reports the discrepancies between reference empirical correlation and generated time series correlation for dimension $d = 4, 10, 20$. We can see that in dimension up to 20, CEGEN obtains a significant improvement compared to all the GANs. Figure 2 illustrates how well the CEGEN generator outperforms the other generators with respect to the FID and QVar metrics in the higher dimensions. This is confirmed by statistics on volatility and drift, as well as by envelope statistics described in Table 6 in Appendix. Figure 5 shows that the 20 processes envelopes are well respected by CEGEN. These good global performances encourage us to focus on the conditional generator in the following transfer learning section.

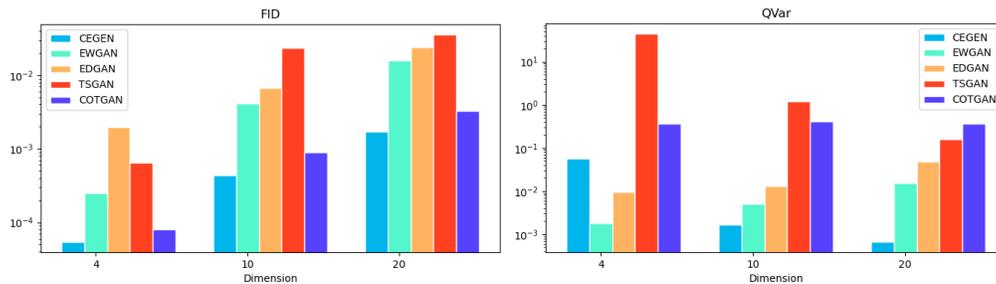


Figure 2: **Exp. B. Left:** Average of Fréchet Inception Distance between distributions at each time step. **Right:** MSE between quadratic variations. (Log scale).

5.5 Transfer learning for small dataset (Exp. C)

Deep generators may need more data than available to be trained effectively. As is done in transfer learning (Torrey and Shavlik, 2010), we propose to start the training with a reasonably wrong model and to finish up the training with the few real data samples. This situation is tested on synthetic data and allows us to track

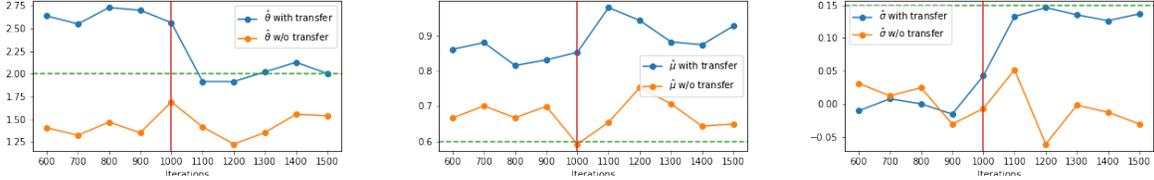


Figure 3: **Exp. C.** Evolution of parameter estimations during training when a transfer occurs at iteration 1000 (red lines). The dashed green lines correspond to the theoretical target value of the real model. The orange lines indicates coefficient estimation of CEGEN only trained on few data and blue lines CEGEN which is first trained on misspecified model then with few real data.

the drift and volatility parameters evolution during the training phase. The reference data are assumed to come from samples simulated from an OU process, while the wrong but reasonable Monte Carlo samples come from a misspecified OU model. The original (resp. misspecified) parameters are $\sigma = 0.15$, $\mu = 0.6$, $\theta = 2.0$ (resp. $\sigma_{MC} = 0.1$, $\mu_{MC} = 0.8$, $\theta_{MC} = 3.0$) and original data sequences include only 60 sequences of 30 dates (5 years of monthly measures). The CEGEN with transfer is compared with a CEGEN only trained with the few available real sequences.

Figure 3 provides the coefficient evolution of both generators during the training process. The transfer iteration time is represented by the red vertical line. Firstly trained with miscalibrated OU model, the transfer learning approach is able to retrieve the parameters when fed with few samples of the target model. The generator only trained with few real samples is unable to estimate correctly the θ and σ coefficients, but exhibits a better estimation of μ . The CEGEN benefiting from the transfer learning takes advantage of the initial training phase and provides an overall better estimation.

This framework is a way to update an existing model with the help of incoming real data. In this situation, the training would start with samples simulated from the consistent model and end with real world inputs.

5.6 Experiments on real-world datasets (Exp. D)

Finally, we test the CEGEN algorithm on various real data with heterogeneous time series. In Table 3, we compare CEGEN performances with the help of FID, QVar and Corr. Our model outperforms GANs or are close in term of FID and QVar for each real times series, and captures well the correlation structure of the majority of the signals. However, some QVar from TSGAN or COTGAN are lower than CEGEN despite their generated trajectories are significantly smoother than real data. To better evaluate the fidelity of the generation we need to consider *post-hoc* metrics. Table 4 reports discriminative and predictive scores for all models (except EWGAN where scores can be found in Appendix), the lower the better. Our generator almost consistently generates higher-quality time series in comparison to the benchmark. On Electric Load data, COTGAN is better able to capture seasonality of the times series, but generates too smooth trajectories. In the opposite CEGEN proposes more faithful times series in term of noise, but struggles to fool the classifier.

Data	CEGEN			EDGAN		
	FID	QVar	Corr	FID	QVar	Corr
Spot prices (d=2)	1.38e-04	2.09e+00	2.10e-02	3.11e-03	2.18e+00	4.12e-02
Stocks (d=6)	1.04e-04	2.10e+01	2.33e-03	7.93e-03	2.43e+01	9.78e-03
Electric Load (d=12)	6.47e-03	4.30e+00	1.27e-03	4.62e-02	1.27e+00	1.56e-03
Jena climate (d=15)	1.10e-03	7.18e+00	1.75e-02	4.39e-02	7.73e+00	1.46e-01
Data	TSGAN			COTGAN		
	FID	QVar	Corr	FID	QVar	Corr
Spot prices (d=2)	2.12e-04	9.00e-02	4.45e-02	1.09e-04	8.25e-01	4.15e-02
Stocks (d=6)	3.46e-03	2.19e+01	2.76e-01	1.49e-04	1.86e+01	1.62e-03
Electric Load (d=12)	5.12e-03	9.18e-01	1.87e-03	4.10e-01	3.45e+00	6.27e-01
Jena climate (d=15)	4.07e-03	8.49e+01	1.89e-02	4.48e-03	7.90e+00	2.34e-02

Table 3: **Exp. D.** Accuracy evaluations for generations on real world time series (the lower, the better).

Data	CEGEN		EDGAN		TSGAN		COTGAN	
	Disc	Pred	Disc	Pred	Disc	Pred	Disc	Pred
Spot prices (d=2)	.014	.049	.137	.049	.066	.055	.033	.049
Stocks (d=6)	.079	.040	.429	.041	.159	.041	.116	.041
Electric Load (d=12)	.433	.028	.495	.046	.407	.032	.277	.022
Jena climate (d=15)	.140	.032	.483	.035	.179	.032	.227	.042

Table 4: **Exp. D.** Discriminative and Predictive scores on real time series (the lower, the better).

Conclusion

We introduced three generative methods for times series, relying on a Deep Euler representation and Wasserstein distances. Two generative method EWGAN and EDGAN demonstrate an accuracy similar to state-of-the-art GAN generators and show better performance for capturing temporal dynamic metrics of the time series. The third method CEGEN is based on a loss metric computed on the conditional distributions of the time series. We prove that minimizing this loss ensures a proper estimation of the drift and volatility coefficients of underlying Itô processes. Our experiments on synthetic and real-world datasets demonstrate that CEGEN outperforms the other generators marginal and temporal dynamics metrics. CEGEN is able to capture correlation structures in high dimensions and is robust when combined with transfer learning on sparse datasets. Transfer learning tests show how this type of method can rely on a proven simulation model without replacing it completely. In further work, we plan to consider more specialized neural networks architectures for time series, extend our results to more general Lévy processes, which may include jumps, and consider not Gaussian noise.

Broader impact

Generative methods for time series may be involved in industries using stochastic control and stochastic simulation methods making them of particular interest for the financial industry, for utilities and energy companies. When applied within a decision-making process, generative methods has to be used carefully as a failure during learning phase may lead to damageable consequences. In this situation, the outputs of the generators should not be left free, as this could lead to erratic optimal controls. Contrarily to the existing approaches which applies GANs and embedding to generate any kind of time series, we impose an Euler structure and we restrain ourselves within the (sufficiently) large class of Itô processes. One of the proposed algorithms manages to get good behavior for synthetic as well as for real data. Moreover, mathematical proofs gives an error estimate of the underlying process parameters for a given loss level.

References

- M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- R. Bhatia, T. Jain, and Y. Lim. On the bures–wasserstein distance between positive definite matrices. *Expositiones Mathematicae*, 37(2):165–191, 2019.
- H. Buehler, L. Gonon, J. Teichmann, and B. Wood. Deep hedging. *Quantitative Finance*, 19(8):1271–1291, 2019.
- H. Buehler, B. Horvath, T. Lyons, I. Perez Arribas, and B. Wood. A data-driven market simulator for small data environments. *Available at SSRN 3632431*, 2020.
- Y. Chen, Y. Wang, D. Kirschen, and B. Zhang. Model-free renewable scenario generation using generative adversarial networks. *IEEE Transactions on Power Systems*, 33(3):3265–3275, 2018.
- A. Clark, J. Donahue, and K. Simonyan. Adversarial video generation on complex datasets. *arXiv preprint arXiv:1907.06571*, 2019a.
- A. Clark, J. Donahue, and K. Simonyan. Efficient video generation on complex datasets. 2019b.
- M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26:2292–2300, 2013.
- C. Donahue, J. McAuley, and M. Puckette. Adversarial audio synthesis. *arXiv preprint arXiv:1802.04208*, 2018.
- B. Efron. *The jackknife, the bootstrap and other resampling plans*. SIAM, 1982.
- C. Esteban, S. L. Hyland, and G. Rätsch. Real-valued (medical) time series generation with recurrent conditional gans. *arXiv preprint arXiv:1706.02633*, 2017.
- S. Fecamp, J. Mikael, and X. Warin. Deep learning for discrete-time hedging in incomplete markets. *Journal of Computational Finance*, 2020.
- A. Fermanian. Embedding and learning with signatures. *arXiv preprint arXiv:1911.13211*, 2019.
- M. Gelbrich. On a formula for the l2 wasserstein metric between measures on euclidean and hilbert spaces. *Mathematische Nachrichten*, 147(1):185–203, 1990.
- A. Genevay, G. Peyré, and M. Cuturi. Learning generative models with sinkhorn divergences. In *International Conference on Artificial Intelligence and Statistics*, pages 1608–1617, 2018.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville. Improved training of wasserstein gans. *arXiv preprint arXiv:1704.00028*, 2017.
- M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *arXiv preprint arXiv:1706.08500*, 2017.
- N. Ikeda and S. Watanabe. *Stochastic differential equations and diffusion processes*. Elsevier, 2014.
- P. Jorion. *Value at risk*. McGraw-Hill Professional Publishing, 2000.
- C. F. Kelliher and L. S. Mahoney. Using monte carlo simulation to improve long-term investment decisions. *The Appraisal Journal*, 68(1):44, 2000.
- D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

- Y. Li, K. Swersky, and R. Zemel. Generative moment matching networks. In *International Conference on Machine Learning*, pages 1718–1727, 2015.
- X. Lyu, M. Hueser, S. L. Hyland, G. Zerveas, and G. Raetsch. Improving clinical predictions through unsupervised time series representation learning. *arXiv preprint arXiv:1812.00490*, 2018.
- L. Malago, L. Montrucchio, and G. Pistone. Wasserstein riemannian geometry of positive definite matrices. *arXiv preprint arXiv:1801.09269*, 2018.
- O. Mogren. C-rnn-gan: Continuous recurrent neural networks with adversarial training. *arXiv preprint arXiv:1611.09904*, 2016.
- S. L. Mullen and D. P. Baumhefner. Monte carlo simulations of explosive cyclogenesis. *Monthly weather review*, 122(7):1548–1567, 1994.
- B. Muzellec and M. Cuturi. Generalizing point embeddings using the wasserstein space of elliptical distributions. *arXiv preprint arXiv:1805.07594*, 2018.
- A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- H. Pham. *Continuous-time stochastic control and optimization with financial applications*, volume 61. Springer Science & Business Media, 2009.
- M. Saito, E. Matsumoto, and S. Saito. Temporal generative adversarial nets with singular value clipping. In *Proceedings of the IEEE international conference on computer vision*, pages 2830–2839, 2017.
- M. Sorge. Stress-testing financial systems: an overview of current methodologies. 2004.
- D. Spehner, F. Illuminati, M. Orszag, and W. Roga. Geometric measures of quantum correlations with bures and hellinger distances. *Lectures on General Quantum Correlations and their Applications*, page 105, 2017.
- S. Steinerberger. Wasserstein distance, fourier series and applications. *arXiv preprint arXiv:1803.08011*, 2018.
- L. Torrey and J. Shavlik. Transfer learning. In *Handbook of research on machine learning applications and trends: algorithms, methods, and techniques*, pages 242–264. IGI global, 2010.
- S. Tulyakov, M.-Y. Liu, X. Yang, and J. Kautz. Mocogan: Decomposing motion and content for video generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1526–1535, 2018.
- C. Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- C. Vondrick, H. Pirsiavash, and A. Torralba. Generating videos with scene dynamics. *Advances in neural information processing systems*, 29:613–621, 2016.
- M. Wiese, R. Knobloch, R. Korn, and P. Kretschmer. Quant gans: Deep generation of financial time series. *Quantitative Finance*, 20(9):1419–1440, 2020.
- J. Xu, X. Ren, J. Lin, and X. Sun. Dp-gan: diversity-promoting generative adversarial network for generating informative and diversified text. *arXiv preprint arXiv:1802.01345*, 2018.
- T. Xu, L. K. Wenliang, M. Munn, and B. Acciaio. Cot-gan: Generating sequential data via causal optimal transport. *arXiv preprint arXiv:2006.08571*, 2020.
- J. Yoon, D. Jarrett, and M. van der Schaar. Time-series generative adversarial networks. 2019.
- Y. Zhang, Z. Gan, K. Fan, Z. Chen, R. Henao, D. Shen, and L. Carin. Adversarial feature matching for text generation. In *International Conference on Machine Learning*, pages 4006–4015. PMLR, 2017.

A Material of Section 4.2

A.1 Counterexample

Consider a timegrid $\{t_0, t_1, t_2\}$ with $\Delta t = t_{i+1} - t_i$. Let $b_X : (t, x) \rightarrow -2x/\Delta t$, $b_Y : (t, x) \rightarrow 0$ and $\sigma_X(t, x) = \sigma_Y(t, x) = 1$ for all $x \in \mathbb{R}$. The b 's and σ 's are Lipschitz in the second coordinate:

$$\begin{aligned}\|b_X(t, x) - b_X(t, y)\|_2 &\leq K\|x - y\|_2 \\ \|b_Y(t, x) - b_Y(t, y)\|_2 &\leq K\|x - y\|_2\end{aligned}$$

For the uni-dimensional case, we have:

$$\begin{aligned}X_{t_0} &= 0, & X_{t_1} &= \varepsilon_{t_1}^X, & X_{t_2} &= -X_{t_1} + \varepsilon_{t_2}^X \\ Y_{t_0} &= 0, & Y_{t_1} &= \varepsilon_{t_1}^Y, & Y_{t_2} &= Y_{t_1} + \varepsilon_{t_2}^Y\end{aligned}$$

with $\varepsilon_{t_i}^X, \varepsilon_{t_i}^Y \sim \mathcal{N}(0, \Delta t)$ and being i.i.d. Then, for $i \in \{0, 1, 2\}$, we have $\mathcal{L}(X_{t_i}) = \mathcal{L}(Y_{t_i})$ but $\mathbb{E}[X_{t_2}|X_{t_1} = z] = -z \neq z = \mathbb{E}[Y_{t_2}|Y_{t_1} = z]$

A.2 Motivation - Details

Proposition A.1. *Assume that for all $t_i \in \{t_0, \dots, t_N = T\}$, for all $z \in \mathbf{R}^d$, $X_{t_i+\Delta t}|X_{t_i} = z$ and $Y_{t_i+\Delta t}|Y_{t_i} = z$ are identically distributed and that $\sigma_X(t_i, z)\sigma_X(t_i, z)^T$ (resp. $\sigma_Y(t_i, z)\sigma_Y(t_i, z)^T$) are positive semi-definite. Then $b_X(t_i, z) = b_Y(t_i, z)$ and $\sigma_X(t_i, z) = \sigma_Y(t_i, z)$.*

Proof. Let $t_i \in \{t_0, \dots, t_n = T\}$. For $z \in \mathbb{R}^d$. We have,

$$\begin{aligned}X_{t_i+\Delta t}|(X_{t_i} = z) &\sim \mathcal{N}(z + b_X(t_i, z)\Delta t, \sigma_X^2(t_i, z)\Delta t) \\ Y_{t_i+\Delta t}|(Y_{t_i} = z) &\sim \mathcal{N}(z + b_Y(t_i, z)\Delta t, \sigma_Y^2(t_i, z)\Delta t)\end{aligned}$$

then, $b_Y(t_i, z) = b_Y(t_i, z)$ and $\sigma_Y(t_i, z)\sigma_Y(t_i, z)^T = \sigma_X(t_i, z)\sigma_X(t_i, z)^T$ for $z \in \mathbb{R}^d$. Matrix $\sigma_X(t_i, z)\sigma_X(t_i, z)^T$ being PSD, has a unique square root which is $\sigma_X(t_i, z)$. The same goes for $\sigma_Y(t_i, z)$. So, $\sigma_Y(t_i, z) = \sigma_X(t_i, z)$. \square

A.3 Proof of Proposition 4.1

Assume that $\sigma_X^2(t_i, \cdot)$, $\sigma_Y^2(t_i, \cdot)$ are strictly positive and, together with $b_X(t_i, \cdot)$ and $b_Y(t_i, \cdot)$, K -Lipschitz in their second coordinate. For $t_i \in \mathcal{T}$, let $(I_k)_k$ be a regular partition covering $\text{Supp}(X_{t_i}) \cup \text{Supp}(Y_{t_i})$ with mesh size Δx . Let $\varepsilon > 0$.

The Itô process X follows the dynamics:

$$X_{t+\Delta t} = X_t + b_X(t, X_t)\Delta t + \sigma_X(t, X_t)\mathcal{N}(0, \Delta t)$$

Thus, for all $z \in \mathbb{R}^d$, $X_{t_i+\Delta t}|(X_{t_i} = z) \sim \mathcal{N}(z + b_X(t_i, z)\Delta t, \sigma_X^2(t_i, z)\Delta t)$ and the same goes for process Y .

Let $I_k = [a_k^1, a_{k+1}^1] \times [a_k^d, a_{k+1}^d]$. Suppose that for all $t_i \in \{t_0, \dots, t_N = T\}$, we have

$$\mathcal{W}_2^2(\mathcal{L}(X_{t_{i+1}}|(X_{t_i} \in I_k)), \mathcal{L}(Y_{t_{i+1}}|(Y_{t_i} \in I_k)) \leq \varepsilon \quad (7)$$

then by definition of \mathcal{W}_2^2 , $\exists z_1, z_2 \in I_k$, such that:

$$\begin{aligned}\|z_1 + b_X(t_i, z_1)\Delta t - z_2 - b_Y(t_i, z_2)\Delta t\|_2^2 \\ + \mathcal{B}^2(\sigma_X^2(t_i, z_1)\Delta t, \sigma_Y^2(t_i, z_2)\Delta t) \leq \varepsilon\end{aligned} \quad (8)$$

By standard norm inequalities and Lipschitz properties of $b_X(t_i, \cdot)$ and $b_Y(t_i, \cdot)$, we bound the squared distance of drifts for all $z \in I_k$ and with mesh grid $\Delta x = \|a_{k+1} - a_k\|_2$.

$$\|(b_X(t_i, z_1) - b_Y(t_i, z_2))\Delta t + (z_1 - z_2)\|_2^2 \leq \varepsilon$$

$$\begin{aligned}
\|b_X(t_i, z_1) - b_Y(t_i, z_2)\|_2 &\leq \frac{\sqrt{\varepsilon} + \Delta x}{\Delta t} \\
\|b_Y(t_i, z_2) - b_X(t_i, z)\|_2 &= \|b_X(t_i, z_1) - b_X(t_i, z)\|_2 \\
&\leq \frac{\sqrt{\varepsilon} + \Delta x}{\Delta t} \\
\|b_Y(t_i, z_2) - b_X(t_i, z)\|_2 &\leq \frac{\sqrt{\varepsilon} + \Delta x}{\Delta t} + K\Delta x \\
\|b_Y(t_i, z) - b_X(t_i, z)\|_2 &\leq \frac{\sqrt{\varepsilon} + \Delta x}{\Delta t} + 2K\Delta x
\end{aligned}$$

We recall that the Bures metrics (Bhatia et al., 2019; Malago et al., 2018) between positive definite matrices A and B is defined by

$$\mathcal{B}^2(A, B) = \text{Tr}(A) + \text{Tr}(B) - 2\text{Tr}(A^{\frac{1}{2}}BA^{\frac{1}{2}})^{1/2}$$

For volatility bound, from Equation (8) we have, In the $\mathbf{d}=\mathbf{1}$ case, this implies $\|\sigma_X(t_i, z_1)\sqrt{\Delta t} - \sigma_Y(t_i, z_2)\sqrt{\Delta t}\|_2^2 \leq \varepsilon$ which leads to : for all $z \in I_k$,

$$\|\sigma_X(t_i, z) - \sigma_Y(t_i, z)\|_2 \leq \sqrt{\frac{\varepsilon}{\Delta t}} + 2K\Delta x$$

For $\mathbf{d}>\mathbf{1}$, let's denote \mathcal{H} the Hellinger distance between positive density matrices:

$$\mathcal{H}(A, B) = \|A^{\frac{1}{2}} - B^{\frac{1}{2}}\|_2 \quad (9)$$

For two density matrices A and B , from (Spehner et al., 2017) (Equation 74) we have $\mathcal{H}(A, B) < \sqrt{2}\mathcal{B}(A, B)$. Following the trace assumption, $\text{Tr}(\sigma_X^2(t_i, z_1)) = \text{Tr}(\sigma_Y^2(t_i, z_2)) = \alpha$ and then,

$$\mathcal{H}\left(\frac{\sigma_X^2(t_i, z_1)}{\alpha}\Delta t, \frac{\sigma_Y^2(t_i, z_2)}{\alpha}\Delta t\right) \leq \sqrt{2\varepsilon} \quad (10)$$

Thus we get,

$$\|\sigma_X(t_i, z_1) - \sigma_Y(t_i, z_2)\|_2 \leq \sqrt{\frac{2\alpha\varepsilon}{\Delta t}}. \quad (11)$$

In particular, using the K -Lipschitz property of the volatility functions, we obtain: for all $z \in I_k$,

$$\|\sigma_X(t_i, z) - \sigma_Y(t_i, z)\|_2 \leq \sqrt{\frac{2\alpha\varepsilon}{\Delta t}} + 2K\Delta x \quad (12)$$

Remark A.2. *The Proposition above extends to the Wasserstein-2 loss if we conditionate from points instead of conditioning from intervals. Indeed, suppose for all $t_i \in \{t_0, \dots, t_N = T\}$, we have*

$$\mathcal{W}_2^2(\mathcal{L}(X_{t_{i+1}}|(X_{t_i} = z_1), \mathcal{L}(Y_{t_{i+1}}|(Y_{t_i} = z_2)) \leq \varepsilon \quad (13)$$

then, as we have :

$$\begin{aligned}
X_{t_i+\Delta t}|(X_{t_i} = z_1) &\sim \mathcal{N}(z_1 + b_X(t_i, z_1)\Delta t, \sigma_X^2(t_i, z_1)\Delta t) \\
Y_{t_i+\Delta t}|(Y_{t_i} = z_2) &\sim \mathcal{N}(z_2 + b_Y(t_i, z_2)\Delta t, \sigma_Y^2(t_i, z_2)\Delta t)
\end{aligned}$$

We can use the closed form of the Gaussian expression of Wasserstein-2:

$$\begin{aligned}
&\|z_1 + b_X(t_i, z_1)\Delta t - z_2 - b_Y(t_i, z_2)\Delta t\|_2^2 \\
&\quad + \mathcal{B}^2(\sigma_X^2(t_i, z_1)\Delta t, \sigma_Y^2(t_i, z_2)\Delta t) \leq \varepsilon
\end{aligned} \quad (14)$$

Similar results as proof 4.1 follow.

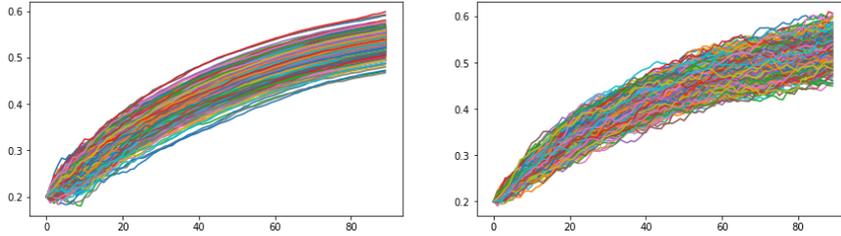


Figure 4: **Left:** Unsatisfactory generations from GAN **Right:** Reference (Ornstein-Uhlenbeck)

B Additional numerical results

B.1 Experiment A - Unidimensional case with synthetic data (Exp. A)

Quantitative evaluations provided by Table 5 for each models highlight how CEGEN and TSGAN stand out in term of average moments accuracy for a 1-dimensional Black-Scholes. We can see that the benchmark TSGAN and our model CEGEN are faithful to synthetic trajectories and outperform both Euler GANs. However, CEGEN is also able to capture the temporal dynamics of both processes, as QVar metrics reports.

Black-Scholes				
metrics	CEGEN	EWGAN	EDGAN	TSGAN
q05	9.95e-05	2.55e-02	3.49e-03	2.61e-06
Avg	8.01e-07	4.04e-02	2.20e-05	9.44e-07
q95	7.85e-06	6.02e-02	2.85e-03	2.47e-05
QVar	4.54e-04	7.30e-02	4.12e-02	2.38e+00
Ornstein-Uhlenbeck				
metrics	CEGEN	EWGAN	EDGAN	TSGAN
q05	4.89e-04	3.98e-03	7.30e-02	2.96e-06
Avg	2.27e-07	2.39e-05	4.47e-02	2.25e-06
q95	8.55e-04	3.98e-03	2.38e-02	6.49e-06
QVar	4.59e-03	2.51e+00	1.67e-03	1.04e+00

Table 5: **Exp. A** Mean squared error (MSE) between reference samples and generated time series on marginal metrics.

B.2 Experiment B - Multidimensional case with synthetic data (Exp. B)

Figure 5 reports envelope of samples from CEGEN model (orange) on a 20-dimensional Black-Scholes (blue). Full lines give the marginal averages over time, and dash lines give average 5% and 95% quantiles respectively. Our generator is still able to retrieve faithfully the average moments in high dimension. This is confirmed with quantitative evaluation provided by Table 6 where the benchmark TSGAN and CEGEN stand out compared to Euler GANs.

Empirical correlation matrices are also retrieved by CEGEN up to dimension 20. In Figure 6, we represent the reference empirical correlation alongside the one coming from generated samples in both correlated and independent case. The term by term mean squared error of correlation matrices (the more black, the better) confirms that CEGEN generations are highly realistic.

B.3 Experiment C - Transfer Learning (Exp. C)

We provide in Table 7 the empirical coefficients of the Ornstein-Uhlenbeck we try to generate. The CEGEN algorithm benefiting of transfer learning gives the closer estimation to the real parameters. However, the CEGEN only trained on the few available samples proposes a better estimation of μ term. A deeper analysis would be welcome and is the subject of future work.

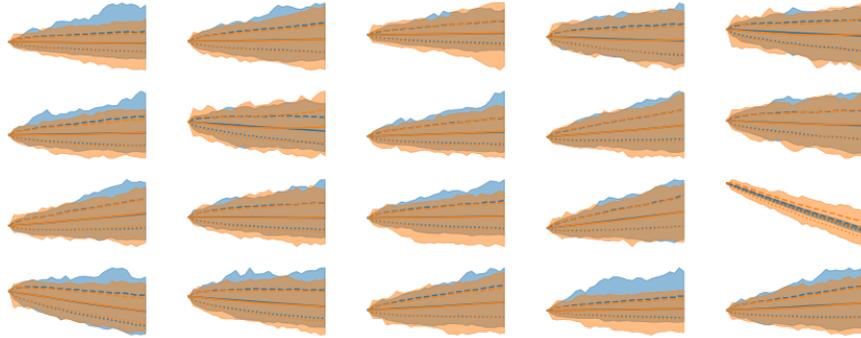


Figure 5: **Exp. B.** Samples from CEGEN model on 20-dimensional Black-Scholes model-based generation. Envelope of both CEGEN (orange) and Reference (blue) time series.

	CEGEN	EWGAN	EDGAN	TSGAN
Dimension = 4				
q05	7.58e-05	2.24e-05	1.39e-03	3.25e-05
Mean	6.02e-06	1.34e-05	1.12e-05	4.21e-06
q95	4.46e-05	1.38e-04	1.23e-03	3.31e-05
Dimension = 10				
q05	1.63e-04	5.23e-04	1.86e-03	2.00e-03
Mean	6.59e-06	2.60e-05	1.27e-05	5.12e-04
q95	2.80e-04	9.17e-04	1.91e-03	6.25e-03
Dimension = 20				
q05	1.16e-04	9.51e-04	3.70e-03	1.19e-04
Mean	1.19e-05	4.88e-05	3.71e-05	1.62e-05
q95	4.35e-04	1.88e-03	2.38e-03	4.14e-04

Table 6: **Exp. B.** Mean squared error (MSE) between reference and generated envelope statistics on a BS case

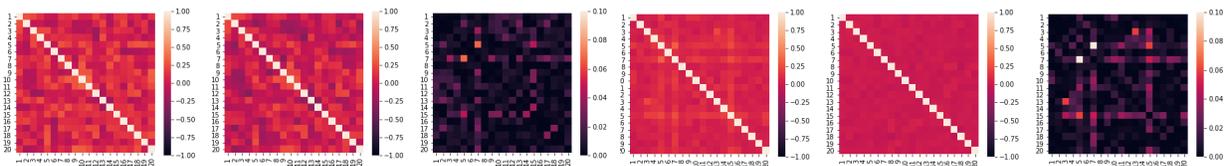


Figure 6: **Exp. B** Mean over time of empirical correlation matrices illustrates a diagonal covariance matrix case (independence case). First heatmap is generated samples from CEGEN (left), second is Monte Carlo ones (right), the target, the third heatmap (mostly black) represents the mean squared error of the two correlation matrices.

	Misspecified model	CEGEN w/o transfer	CEGEN with transfer
θ (2.00)	3.00	1.54	1.78
μ (0.60)	0.80	0.65	1.01
σ (0.15)	0.10	-0.03	0.15

Table 7: **Exp. C.** Empirical Ornstein-Uhlenbeck coefficient estimations according to misspecified samples, CEGEN trained on few original data and CEGEN with transfer learning.

B.4 Experiment D - Real data and other benchmarks (Exp. D)

Table 8 reports discriminative and predictive performances of EWGAN, the conditional recurrent RCGAN (Esteban et al., 2017) and an unconditional MMD with Gaussian kernel GMMN (Li et al., 2015).

Data	EWGAN		RCGAN		GMMN	
	Disc	Pred	Disc	Pred	Disc	Pred
Spot prices (d=2)	.225	.050	.427	.809	.137	.671
Stocks (d=6)	.238	.042	.287	.616	.499	.626
Electric Load (d=12)	.410	.029	.495	.581	.499	.566
Jena climate (d=15)	.479	.034	.499	.651	.295	.634

Table 8: **Exp. D.** Discriminative and Predictive scores on real time series (the lower, the better).

C Algorithms details

In this section, we detail the pseudo algorithms of both Euler GANS (EWGAN and EDGAN). We also explain more deeply the conditional loss computation, as well as some tested variants.

Algorithm 2 Euler Wasserstein-1 Generative Adversarial Networks (EWGAN).

Input: θ_0, φ_0 randomly chosen, α, β , learning rates,
 K number of iterations, M batch size, n_c critic iterations, c clipping value, $(X^{(i)})_{i=1..M}$ real data
Output: θ, φ
 $\theta \leftarrow \theta_0, \varphi \leftarrow \varphi_0$
while NOT CONVERGE **do**
 for $j = 1..n_{\text{critic}}$ **do**
 $x \leftarrow M$ samples with $X^{(i)} = (X_{t_1}^{(i)}, \dots, X_{t_N}^{(i)})_{i=1..M}$
 $z \leftarrow M$ samples iid gaussian noise
 $y^\theta \leftarrow M$ generation from Euler scheme and $g_\theta(z)$
 $\varphi \leftarrow \varphi + \alpha \text{Adam}(\nabla_\varphi(\mathbb{E}[d_\varphi(x)] - \mathbb{E}[d_\varphi(y^\theta)]), \alpha)$
 $\varphi \leftarrow$ gradient penalty $(\varphi, 10)$
 end for
 $x \leftarrow M$ samples with $X^{(i)} = (X_{t_1}^{(i)}, \dots, X_{t_N}^{(i)})_{i=1..M}$
 $z \leftarrow M$ samples iid Gaussian noise
 $Y^\theta \leftarrow M$ generation from Euler scheme and $g_\theta(z)$
 $\theta \leftarrow \theta - \beta \text{Adam}(\nabla_\theta \mathbb{E}[d_\varphi(Y^\theta)], \beta)$
end while

For given loss ℓ (Bures-Wasserstein (6) in the paper), we compute the conditional loss by extracting the elements at a certain date such that the previous state belongs to an ensemble I . We propose below two ways to do it.

The first approach consists in sorting each dimension at each time step in order to get K quantiles $(a_k)_{k=1..K}$ of both X_{t_i} and Y_{t_i} . At date t_i for $i \in \{1, \dots, T\}$ and for each dimension $d \in \{1, \dots, D\}$, for a given batch of samples, $\mathcal{L}(X_{t_{i+1}} | X_{t_i} \in I)$ is approximated by selecting only the realizations $x_{t_{i+1}}^d$ such that the previous state $x_{t_i}^d$ belongs to the interval $I_k^d = [a_k^d, a_{k+1}^d]$. The losses ℓ_k^d between the two conditional distributions are then summed up over all dimensions and subdivisions. To take into account the disjoint support case (for instance samples $x_{t_i}^d \in [-1, 0[$ and $y_{t_i}^d \in]0, 1]$), we penalize by the distance separating the supports. See Algorithm 4 for further details.

Another approach is to compute the partitions of $\text{Supp}(X_t^d)$ before the generator training phase. We use T Kmeans to compute the centers of K clusters at each time step. Then, during the generator training we compute the loss ℓ between $X_{t_{i+1}}$ and $Y_{t_{i+1}}$ such that their respective previous states belong to the same cluster k . This method has the advantage that the generated samples share the same support as the real data one. We use the conditional loss by disjoint quantiles in our experiments, because the algorithm runs faster and gives better empirical results.

Algorithm 3 Euler Dual Generative Adversarial Networks (EDGAN).

Input: θ_0, φ_0 randomly chosen, α, β, γ learning rates,
 K number of iterations, M batch size, n_c critic iterations, c clipping value, $(X^{(i)})_{i=1..M}$ real data
Output: θ, φ, ψ $\theta \leftarrow \theta_0, \varphi \leftarrow \varphi_0, \psi \leftarrow \psi_0$;;
while NOT CONVERGE **do**
 for $j = 1..n_{\text{critic}}$ **do**
 $x \leftarrow M$ samples of $X^{(i)} = (X_{t_1}^{(i)}, \dots, X_{t_N}^{(i)})_{i=1..M}$
 $z \leftarrow M$ samples iid Gaussian noise
 $y^\theta \leftarrow M$ generations from Euler scheme and $g_\theta(z)$
 $\varphi \leftarrow \varphi + \alpha \text{Adam}(\nabla_\varphi(\mathbb{E}[d_\varphi(x)] - \mathbb{E}[d_\varphi(y^\theta)])); \alpha$
 $\varphi \leftarrow$ gradient penalty ($\varphi, 10$)
 for $t = t_1..t_N$ **do**
 $\psi \leftarrow \psi + \gamma \text{Adam}(\nabla_\psi(\mathbb{E}[d_\psi(x_t)] - \mathbb{E}[d_\psi(y^\theta)_t])); \gamma$
 $\psi \leftarrow$ gradient penalty ($\psi, 10$)
 end for
 end for
 $x \leftarrow M$ samples with $X^{(i)} = (X_{t_1}^{(i)}, \dots, X_{t_N}^{(i)})_{i=1..M}$
 $z \leftarrow M$ samples iid Gaussian noise
 $Y^\theta \leftarrow M$ generations from Euler scheme and $g_\theta(z)$
 $\theta \leftarrow \theta - \beta \text{Adam}(\nabla_\theta \mathbb{E}[d_\varphi(Y^\theta)], \beta)$
end while

Algorithm 4 Conditional Loss by disjoint quantiles.

Input: processes of length T $X = (X^1, \dots, X^D), Y = (Y^1, \dots, Y^D), \lambda$
for $t = 1..T$ **do**
 for $d = 0..D$ **do**
 $I_{K,x}^d \leftarrow K$ subdivisions of $\text{Supp}(X_t^d)$;
 $I_{K,y}^d \leftarrow K$ subdivisions of $\text{Supp}(Y_t^d)$;
 for $k = 0..K$ **do**
 if $\text{Supp}(X_t^d) \cup \text{Supp}(Y_t^d) \neq \emptyset$ **then**
 $\ell_{t+1,k}^d \leftarrow \mathcal{W}_2^2(\mathcal{L}(X_{t+1}|X_t^d \in I_{K,x}^d), \mathcal{L}(Y_{t+1}|Y_t^d \in I_{K,y}^d))$
 else
 $\ell_{t+1,k}^d \leftarrow \lambda |\mathbb{E}[X_t^d] - \mathbb{E}[Y_t^d]|$
 end if
 end for
 end for
end for
 $\ell = \sum_{t=1}^T \sum_{d=1}^D \sum_{k=1}^K \ell_{t+1,k}^d$
Output: ℓ

D Models and hyperparameters

We use tensorflow to implement neural networks. The networks architecture is composed of 3-layers of 4 times the data dimension neurons each (for stocks $4 \times 6 = 24$ neurons). Euler generator networks are feed-forward, as we want to be Markovian, while benchmarks TSGAN and COTGAN architecture benefit of recurrent networks (GRU, LSTM). Code of TSGAN is available online ([link](#)), as well as the code of COTGAN ([link](#)). Other details are precised in Table 9. Real dataset are normalized with MinMax scaler $((x - \min) / (\max - \min))$ and the first date always starts at spot $X_0 = 0.2$.

Settings of neural networks	
T (ndates)	30
White noise dim	$(T \times d)$
Optimizer	Adam
Nb iterations	5000
Batch size	300
Learning rates	1.10^{-3}

Table 9: Neural network hyper-parameters

To compute the discriminative and predictive scores, we use the same network architecture and parameters as (Yoon et al., 2019) (actually we use their code). The neural networks are 2-layer LSTMs with hidden dimensions 4 times the size of the input features, and use tanh as the activation function and sigmoid as the output layer activation function (such that output belongs to $[0,1]$).

The training is done on 12 i7-9750H processors at 2.60 GHz.

E Results variation

Table 10 illustrates the variation between three different trainings of each generators for stocks data. We recall that the discriminative and predictive score are obtained by training 10 LSTM networks and averaging their scores (they thus include some additional variation).

Stocks data	Discriminative	Predictive
EWGAN	.417(\pm .041)	.041(\pm .001)
EDGAN	.444(\pm .146)	.041(\pm .000)
CEGEN	.077(\pm .015)	.040(\pm .000)
TSGAN	.168(\pm .025)	.041(\pm .001)
COTGAN	.094(\pm .022)	.041(\pm .000)

Table 10: Performance variations of each generator for stocks data on discriminative and predictive scores, for three different trainings.

F Data

Dataset	Sequences	Seq. length	Dim.
Price	52608	30	2
Stocks	3600	30	6
Electric Load*	50000	30	13
Jena Climate*	50000	30	15

Table 11: Data description.

Table 11 reports the number of observations of each dataset, the sequence length chosen in our experiments, as well as their dimension. All datasets are available online, and can be downloaded from: RTE for electric load and price ([link](#)), Keras for Jena climate ([link](#)). Stocks data source are described in (Yoon et al., 2019). *For Electric Load and Jena Climate we take only the first 50000 observations.